

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 07-244597

(43)Date of publication of application : 19.09.1995

(51)Int.Cl.

G06F 11/20  
G06F 3/06  
G06F 13/00

(21)Application number : 06-301146

(71)Applicant : INTERNATL BUSINESS MACH  
CORP <IBM>

(22)Date of filing : 05.12.1994

(72)Inventor : KERN ROBERT F  
KERN RONALD MAYNARD  
MCBRIDE GREGORY EDWARD  
SHACKELFORD DAVID MICHAEL

(30)Priority

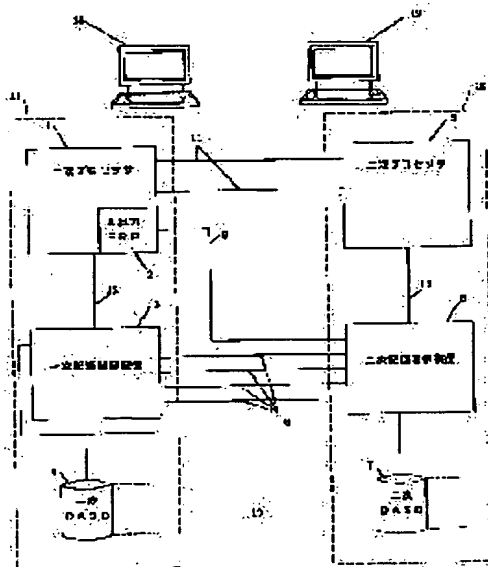
Priority number : 94 199444 Priority date : 22.02.1994 Priority country : US

### (54) METHOD FOR FORMING CONSISTENCY GROUP TO PROVIDE DISASTER RECOVERY FUNCTION, AND ITS RELATED SYSTEM

(57)Abstract:

**PURPOSE:** To provide a remote data shadowing system which provides a real-time disaster recovery function on a storage area base.

**CONSTITUTION:** A write input-output operation is performed in a storage sub-system on the primary side 14 by record update on the primary side 14. A time-stamp is attached to this write input-output operation and the time, order and physical position of the record update are collected in a primary data mover. The primary data mover divides plural sets of record update and their related control information into groups based on prescribed time intervals, adds a prefix header to the record update and thereby forms a self-description



record set. The self-description record set is sent to a remote secondary side 15, and such a consistency group is formed that the record update is ordered to be able to shadow the record update in the sequence that matches the sequence where the write input-output operation was performed on the primary side 14 by the record update.

---

#### LEGAL STATUS

[Date of request for examination] 05.12.1994

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number] 3149325

[Date of registration] 19.01.2001

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 特 許 公 報 (B 2)

(11) 特許番号

特許第3149325号

(P3149325)

(45) 発行日 平成13年3月26日 (2001. 3. 26)

(24) 登録日 平成13年1月19日 (2001. 1. 19)

(51) Int.Cl.<sup>7</sup>

識別記号

F I

G 0 6 F 11/20

3 1 0

G 0 6 F 11/20

3 1 0 C

3/06

3 0 4

3/06

3 0 4 F

12/00

5 3 3

12/00

5 3 3 J

12/16

3 1 0

12/16

3 1 0 M

請求項の数10(全 29 頁)

(21) 出願番号

特願平6-301146

(22) 出願日

平成6年12月5日 (1994. 12. 5)

(65) 公開番号

特開平7-244597

(43) 公開日

平成7年9月19日 (1995. 9. 19)

審査請求日

平成6年12月5日 (1994. 12. 5)

(31) 優先権主張番号

1 9 9 4 4 4

(32) 優先日

平成6年2月22日 (1994. 2. 22)

(33) 優先権主張国

米国 (US)

(73) 特許権者

390009531

インターナショナル・ビジネス・マシー  
ンズ・コーポレーション

INTERNATIONAL BUSI  
NESS MACHINES COR  
PORATION

アメリカ合衆国10504、ニューヨーク州

アーモンク (番地なし)

(72) 発明者

ロバート・フレデリック・カーン

アメリカ合衆国85730 アリゾナ州ツー  
ソン イースト・カレヒコ・ストリート  
8338

(74) 復代理人

100065455

弁理士 山本 仁朗 (外2名)

審査官

多賀 実

最終頁に続く

(54) 【発明の名称】 災害復旧機能を提供するために整合性グループを形成する方法および関連するシステム

(57) 【特許請求の範囲】

【請求項1】 一次データ・ムーバおよびレコード更新を発生する複数のアプリケーションを実行する一次プロセッサを備える一次側と、一次プロセッサと通信可能に接続された二次プロセッサを備える二次側とを含む災害復旧のための遠隔データ・シャドーイング・システムであって、一次プロセッサが一次記憶サブシステムに連結され、一次プロセッサから一次記憶サブシステムに発する書込み入出力操作に従ってレコード更新を記憶するための記憶装置が一次記憶サブシステムに備えられ、一次側が一次側における時間依存操作を同期させるための共通システム・タイマを更に備えており、二次側が順序整合性のある順序でデータ更新のコピーを記憶するための二次記憶サブシステムを備えているシステムにおける順序整合性のある順序でデータ更新をシャドーイングする方

法において、前記方法が、

(a) 一次記憶サブシステムで発生する各書込み入出力操作にタイム・スタンプを付けるステップと、

(b) 一次記憶サブシステムからレコード更新についてのレコード・セット情報を収集するステップと、

(c) データ更新およびレコード・セット情報を一次データ・ムーバに書込んでレコード・セットを生成ステップと、

(d) 二次システムにおける書込み入出力操作の順序を再生成するために二次プロセッサにより使用される自己記述レコード・セットを生成するためレコード・セットの各々にヘッダーを付するステップと、

(e) 所定の時間間隔に従う時間間隔グループで自己記述レコード・セットを二次プロセッサに転送するステップと、

(f) 自己記述レコード・セットの時間間隔グループから整合性グループを形成するステップであって、一次記憶サブシステムに発せられた書込み入出力操作の時間順序に基づいて整合性グループ内においてレコード更新が順序づけされる整合性グループを生成するステップと、

(g) 各整合性グループのレコード更新を整合性ある順序で二次記憶サブシステムにシャドーイングするステップと、

(h) ステップ(e)において受取った各自己記述レコード・セットが完全なものであるかどうかを二次側で判定するステップと、  
を含む方法。

【請求項2】レコード・セットが二次プロセッサに非同期的に転送されることを特徴とする、請求項1に記載の方法。

【請求項3】ステップ(f)が二次側で行われることを特徴とする、請求項1に記載の方法。

【請求項4】自己記述レコード・セットが不完全であると二次側が判定した場合に、欠落データ更新を再送信するよう二次側が一次側に要求することが、ステップ

(e)にさらに含まれることを特徴とする、請求項1に記載の方法。

【請求項5】各時間間隔グループが完全なものであるかどうかを二次側で判定するステップ(i)をさらに含むことを特徴とする、請求項1に記載の方法。

【請求項6】間隔グループが不完全であると二次側が判定した場合に、欠落レコード・セットを再送信するよう二次側が一次側に要求することが、ステップ(i)にさらに含まれることを特徴とする、請求項5に記載の方法。

【請求項7】ステップ(b)が、レコード・セット情報において、各レコード更新が格納されている一次記憶装置上の物理的位置を識別することを特徴とする、請求項1に記載の方法。

【請求項8】ステップ(b)が、レコード・セット情報において、セッション内で一次記憶装置に格納された各レコード更新の順序と更新時間を識別することを特徴とする、請求項7に記載の方法。

【請求項9】ステップ(d)が、接頭部ヘッダにおいて、セッションに関する間隔グループ番号と、そこで参照される各レコード更新のグループ内順序を識別することを特徴とする、請求項1に記載の方法。

【請求項10】レコード更新を生成するアプリケーションを実行する一次側を含み、一次側から離れた位置に二次側を有し、二次側がレコード更新をシャドーイングして、一次側に災害復旧を提供する、リアルタイム・データ・シャドーイングを行う非同期遠隔データ二重化システムにおいて、非同期遠隔データ二重化システムが、一次側の時間依存プロセスを同期させるためのシスプレックス・タイマーと、

アプリケーションを実行し、対応するレコード更新用の書込み入出力操作を出し、一次データ・ムーバをそこに有する、一次側の一次プロセッサと、

各レコード更新ごとに書込み入出力操作を1つずつ受け取る複数の一次記憶制御装置であって、それぞれの一次記憶制御装置書込み入出力操作が一次プロセッサによってシスプレックス・タイマーと同期される、複数の一次記憶制御装置と、

対応する書込み入出力操作に応じて、レコード更新をそこに格納するための複数の一次記憶装置とを含み、

一次データ・ムーバが、各レコード更新ごとに複数の一次記憶制御装置からレコード・セット情報を収集して、所定のグループのレコード・セット情報に接頭部ヘッダを付加し、接頭部ヘッダと所定のレコード・セット情報

グループが自己記述レコード・セットを形成し、各レコード・セット情報が、一次装置アドレス、シリンダ番号およびヘッド番号(CCHH)、レコード更新順序番号、書込み入出力タイプ、検索指数、セクタ番号、およびレコード更新時間を含み、接頭部ヘッダが、総データ長、操作タイム・スタンプ、時間間隔グループ番号、およびレコード読取り時間を含み、

一次側から自己記述レコード・セットを受け取る二次データ・ムーバを有する二次側の二次プロセッサと、二次プロセッサに連結された複数の二次記憶制御装置と、

レコード更新を格納する複数の二次記憶装置とをさらに含み、

二次データ・ムーバが、送信された自己記述レコード・セットが完全なものであるかどうかを判定し、自己記述レコード・セットから整合性グループを形成し、複数の一次記憶装置にレコード更新が書き込まれたときの順序に整合する順序で複数の二次記憶装置に書き込むために各整合性グループから得たレコード更新を複数の二次記憶制御装置に出力する、非同期遠隔データ二重化システム。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は、一般的には、災害復旧技法に関し、より具体的には、直接アクセス記憶装置(DASD)データのリアルタイム遠隔コピーのためのシステムに関する。

【0002】

【従来の技術】データ処理に関連して、効率よくアクセスし、修正および再格納できる大量のデータ(またはレコード)を格納するためには、一般に、データ処理システムが必要である。データ記憶域は、効率よくしかも費用効果の高いデータ記憶を行うために、通常、複数のレベルに、すなわち、階層状態に分かれている。第一のレベル、すなわち、最高レベルのデータ記憶域は、通常はダイナミックまたはスタティック・ランダム・アクセス

・メモリ（DRAMまたはSRAM）である電子メモリを含む。電子メモリは、ナノ秒単位でこのようなバイト数のデータにアクセスすることによってそれぞれの回路に数百万バイトのデータを格納できる、半導体集積回路の形態をとる。アクセスが完全に電子的に行われるため、このような電子メモリは最高速のデータ・アクセスを提供する。

【0003】第二レベルのデータ記憶域は、通常、直接アクセス記憶装置（DASD）を含む。DASD記憶域は、磁気ディスクまたは光ディスクなどを含む可能性があり、これらのディスクは、データのビットを構成する「1」と「0」を表す磁氣的または光学的に変質させたマイクロメートル規模のスポットとしてビット単位のデータをディスク表面に格納する。磁気DASDは、残留磁気材料で被覆した1枚または複数枚のディスクを含む。これらのディスクは、保護環境内に回転式に取り付けられている。各ディスクは、多くの同心トラックまたは間隔が詰まった円に分割されている。データは、各トラックに沿って1ビットずつ順次格納される。ヘッド・ディスク・アセンブリ（HDA）として知られているアクセス機構は、一般に、1つまたは複数の読取り書込みヘッドを含み、ディスクが回転して読取り書込みヘッドを通過する際にディスクの表面にデータを転送したりディスクの表面からデータを転送するためにトラック間を移動できるよう、各DASDに設けられている。通常、ミリ秒単位（電子メモリより低速の単位）でこのようなデータにアクセスすることによって、DASDはギガバイト規模のデータを格納できる。所望のデータ記憶位置にディスクおよびHDAを物理的に配置する必要があるため、DASDに格納されたデータへのアクセス速度は低下する。

【0004】第三またはそれ以下のレベルのデータ記憶域は、テープまたはテープとDASDライブラリを含む。このレベルの記憶域では、必要なデータ記憶媒体を選択して装填するのにロボットが必要になるため、ライブラリ内のデータへのアクセス速度はさらに低下する。利点は、テラバイト規模のデータ記憶など、データ記憶容量が非常に大きい割に費用が低減される点である。テープ記憶装置は、バックアップ目的で 사용되는ことが多い。すなわち、階層の第二レベルに格納されたデータは、磁気テープで安全に保管するために複製される。テープまたはライブラリに格納したデータへのアクセス速度は、現在、秒単位である。

【0005】データの紛失は、企業にとって破滅的なものになる恐れがあるので、多くの企業ではバックアップのデータ・コピーを用意することが必須になっている。一次記憶レベルで紛失したデータの復旧に要する時間も、復旧に関して考慮すべき重要な項目である。テープまたはライブラリでのバックアップ速度の改善策としては、二重コピーが挙げられる。二重コピーの例として

は、追加のDASDにデータが書き込まれるように追加のDASDを用意する方法がある（ミラーリングと呼ばれる場合もある）。この場合、一次DASDで障害が発生しても、データを得るために二次DASDを頼りにすることができる。この方法の欠点は、必要なDASDの数が二倍になることである。

【0006】記憶装置を二重に設ける必要性を克服するもう1つのデータ・バックアップ方法としては、低価格装置の冗長アレイ（RAID）構成にデータを書き込む方法がある。この方法では、多くのDASDにデータを割り当てるようにデータが書き込まれる。1つのDASDで障害が発生しても、残りのデータとエラー修正手順を使用すれば、紛失したデータを復旧できる。現在では、数種類のRAID構成が使用できる。

【0007】一般に、記憶装置または記憶媒体で障害が発生した場合のデータの復旧には、上記のバックアップ方法で十分である。このようなバックアップ方法では、二次データは一次データのミラーになる、つまり、二次データは一次データと同じボリューム通し番号（VOLSER）とDASDアドレスを持つので、これらの方法は装置障害の場合だけ有用である。しかし、ミラーリングされた二次データを使用しても、システム障害の復旧は行えない。このため、地震、火災、爆発、台風など、システム全体またはシステム使用現場を破壊する災害が発生した場合にデータを復旧するには、さらに保護が必要である。つまり、災害復旧のためには、一次データから離れた位置にデータの二次コピーを格納しておく必要がある。災害保護を行う既知の方法としては、毎日または毎週、テープにデータをバックアップする方法がある。この場合、車両でテープを収集し、通常、一次データ位置から数キロメートル離れた安全な保管区域にテープを移送する。このバックアップ計画には、バックアップ・データを検索するのに何日もかかると同時に、数時間分または数日分のデータが失われ、最悪の場合、同じ災害によって保管場所も破壊されてしまう恐れがあるという問題がある。多少改善されたバックアップ方法としては、毎晩、バックアップ場所にデータを転送する方法が考えられる。この方法では、より離れた遠隔地にデータを格納することができる。しかし、この場合も、二重コピー方法と同様、バックアップが連続して行われるわけではないので、次のバックアップまでにデータの一部が失われる可能性がある。このため、一部のユーザにとっては受け入れられないほど、相当なデータ量が失われる恐れがある。

【0008】さらに最近導入されたデータ災害復旧方法としては、遠隔方式だけでなく、連続方式でもデータがバックアップされる、遠隔二重コピーがある。あるホスト・プロセッサから別のホスト・プロセッサへ、またはある記憶制御装置から別の記憶制御装置へ、あるいはこれらを組み合わせた方法で二重データを送信するには、

このプロセスを実現するために相当な量の制御データが必要になる。しかし、オーバーヘッドが高いために、二次側が一次側の処理に遅れないようにする能力が妨げられる可能性があり、このため、災害が発生した場合に二次側が一次側を復旧する能力が脅かされる。

【0009】したがって、最小限の制御データを使用して、一次処理側のデータと一致するデータのリアルタイム更新を行う方法および装置を提供することが望まれている。この方法および装置は、復旧される特定のアプリケーション・データから独立して、つまり、特定のアプリケーション・データをベースとするのではなく、汎用記憶媒体をベースとして機能する。

【0010】

【発明が解決しようとする課題】本発明の目的は、災害復旧のためにDASDデータを二次側にシャドーイング(shadowing)するための改良された設計および方法を提供することにある。

【0011】

【課題を解決するための手段】本発明の第一の実施例によれば、整合性グループを形成するための方法により、遠隔地からの災害復旧機能が提供される。一次プロセッサで実行される1つまたは複数のアプリケーションによって生成されるデータ更新は、一次記憶サブシステムによって受け取られ、この一次記憶サブシステムは、入出力書き込み操作により各データ更新がそこに書き込まれるようにする。一次記憶サブシステムは共通タイマによって同期が取られ、一次プロセッサから離れた位置にある二次システムは、災害復旧のために二次側が使用できるように順序整合性のある順序でデータ更新をシャドーイングする。本方法は、(a)一次記憶サブシステムで行われる各書き込み入出力操作にタイム・スタンプを付けるステップと、(b)各データ更新ごとに一次記憶サブシステムから書き込み入出力操作のレコード・セット情報を収集するステップと、(c)自己記述レコード・セットが一連の書き込み入出力操作を再生成するのに十分なものになるように、データ更新とそれぞれのレコード・セット情報から自己記述レコード・セットを生成するステップと、(d)所定の間隔しきい値に基づいて、自己記述レコード・セットを間隔グループにグループ分けするステップと、(e)最も早い操作タイム・スタンプを持つ自己記述レコード・セットの間隔グループとして、第一の整合性グループを選択するステップとを含み、個々のデータ更新が、一次記憶サブシステム内の入出力書き込み操作の時間順に基づいて、第一の整合性グループ内で順序付けされる。

【0012】本発明の他の実施例では、一次システムは1つまたは複数のアプリケーションを実行する一次プロセッサを有し、このアプリケーションがレコード更新を生成し、一次プロセッサがそれに基づく自己記述レコード・セットを生成する。それぞれの自己記述レコード・

セットは、一次システムから離れた位置にある二次システムに送られ、二次システムは、リアルタイム災害復旧のために、自己記述レコード・セットに基づいて、順序整合性のある順序でレコード更新をシャドーイングする。一次プロセッサは一次記憶サブシステムに接続され、一次記憶サブシステムはレコード更新を受け取って、入出力書き込み操作によって各レコード更新がそこに格納されるようにする。一次プロセッサは、同期を取るためにアプリケーションと一次記憶サブシステムに共通の時間源を提供するためのシスプレックス・クロックを含み、シスプレックス・クロックによって同期が取られる一次データ・ムーバは、各レコード更新ごとに一次データ・ムーバにレコード・セット情報を提供するよう、一次記憶サブシステムに指示する。一次データ・ムーバは、複数のレコード更新とそれに対応する各レコード・セット情報を時間間隔グループにグループ分けし、それに接頭部ヘッダを挿入する。それぞれの時間間隔グループは自己記述レコード・セットを形成する。

【0013】本発明の上記およびその他の目的、特徴、および利点は、添付図面に図示する本発明の実施例に関する以下の詳細な説明から明らかになるだろう。

【0014】

【実施例】一般的なデータ処理システムは、データを計算して操作し、データ機能記憶管理サブシステム/多重仮想記憶システム(DFSMS/MVS)ソフトウェアなどを実行するために、IBMシステム/360またはIBMシステム/370プロセッサなどのホスト・プロセッサの形態をとり、少なくとも1台のIBM 3990記憶制御装置がそれに接続され、その記憶制御装置が、メモリ制御装置と、それに組み込まれた1つまたは複数のタイプのキャッシュ・メモリを含む場合がある。さらに記憶制御装置は、IBM 3380または3390 DASDなどの1群の直接アクセス記憶装置(DASD)に接続されている。ホスト・プロセッサが実質的な計算能力を提供するのに対し、記憶制御装置は、大規模データベースを効率よく転送し、ステージ(stage)/デステージ(destage)し、変換し、全般的にアクセスするのに必要な諸機能を提供する。

【0015】一般的なデータ処理システムの災害復旧保護では、一次DASDに格納した一次データを二次側または遠隔地でバックアップする必要がある。一次側と二次側との距離は、ユーザが受け入れられる危険のレベルによって決まり、数キロメートルから数千キロメートルの範囲が可能である。二次側または遠隔地は、バックアップ・データ・コピーを提供するだけでなく、一次システムが使用不能になった場合に一次システムの処理を引き継ぐのに十分なシステム情報を持っていなければならない。その理由は、主に、単一の記憶制御装置では一次側と二次側に設けた一次および二次DASDストリングの両方にデータを書き込めないためである。むしろ、一

次記憶制御装置に接続された一次DASDストリングには一次データが格納されるのに対し、二次記憶制御装置に接続された二次DASDストリングには二次データが格納されるのである。

【0016】二次側は、一次側から十分離れている必要があるだけでなく、一次データをリアルタイムでバックアップできなくてはならない。二次側は、一次データが更新されたときに、最小限の遅延で一次データをバックアップする必要がある。しかも、二次側は、一次側で実行され、データまたは更新を生成するアプリケーション・プログラム（たとえば、IMSやDB2）を考慮せずに、一次データをバックアップしなければならない。二次側に要求される難しい課題は、二次データの順序が整合していなければならないことである。つまり、二次データは一次データと同じ順序でコピーされなければならない（順序整合性）、これはシステムについてかなり考慮を要する問題である。順序整合性は、それぞれが1つのデータ処理システム内の複数のDASDを制御する記憶制御装置が複数存在するためにさらに複雑になっている。順序整合性がないと、一次データと一致しない二次データが生成され、その結果、災害復旧が崩壊する恐れがある。

【0017】遠隔データの二重化は、同期と非同期の2つの一般的なカテゴリに分けられる。同期遠隔コピーでは、一次データを二次側に送り、一次DASDの入出力操作を終了する（一次ホストにチャネル終了（CE）と装置終了（DE）を出力する）前にこのようなデータの受取りを確認する必要がある。このため、同期コピーでは、二次側の確認を待っている間に一次DASDの入出力応答時間が遅くなる。一次側の入出力応答遅延は、一次システムと二次システムとの距離に比例して長くなる（これは遠隔距離を数十キロメートル規模に制限する要素である）。しかし、同期コピーは、システム・オーバーヘッドを比較的小さくして、順序が整合したデータを二次側に提供する。

【0018】非同期遠隔コピーでは、二次側でデータが確認される前に一次DASDの入出力操作が完了する（一次ホストにチャネル終了（CE）と装置終了（DE）を出力する）ため、一次側のアプリケーション・システムのパフォーマンスが向上する。このため、一次DASDの入出力応答時間は二次側までの距離に依存せず、一次側から数千キロメートル離れた遠隔地に二次側を設けることもできる。しかし、二次側で受け取ったデータの順序が一次側の更新順序と一致なくなる場合が多いので、データの順序整合性を確保するのに必要なシステム・オーバーヘッドが増加する。したがって、一次側で障害が発生すると、一次側と二次側との間で転送中のデータが一部紛失する恐れがある。

#### 【0019】同期データ・シャドーイング

災害復旧のための同期リアルタイム遠隔コピーでは、コ

ピーされた複数のDASDボリュームが1つのセットを形成する必要がある。さらに、このようなセットを形成するには、各セットを構成するこれらのボリューム（VOLSER）とそれに対応する一次側の同等物を識別するために十分な量のシステム情報を二次側に提供する必要がある。重要なのは、二次側が一次側と「二重対（duplex pair）」を形成するので、1つまたは複数のボリュームがこのセットと同期していない、つまり、「二重対の障害」が発生したときに二次側がそれを認識しなければならない点である。代替経路が再試行される間に一次DASDの入出力が遅延するため、非同期遠隔コピーより同期遠隔コピーの方が、接続障害を認識しやすい。一次側は、二次側用の更新を待ち行列に入れる間、一次側が続行できるようにコピーを中止または中断することができ、一次側は、二次側が同期していないことを示すためにこのような更新にマークを付ける。二次側がいつでも災害復旧に使用できるようにするため、二次側が一次側と同期しなくなる原因になりそうな例外条件を認識することが必要である。エラー条件や復旧処置によって、二次側が一次側と一致しなくなってしまう。

【0020】しかし、二次DASDが存在しアクセス可能である状態で、二次側と一次側の接続を維持しても、内容の同期は確保されない。いくつかの理由から、二次側は一次側との同期性を失う場合もある。二重対が形成されたときに二次側は当初同期しておらず、初期データ・コピーが完了したときに同期に達する。一次側が二次側に更新済みデータを書き込めない場合、一次側は二重対を解除することもある。この場合、一次側は、更新アプリケーションが続行できるように、二重対が中断された状況で一次DASDに更新内容を書き込む。このため、二重対が復元されるまで、一次側は露出状態で、つまり、現行の災害保護コピーを使わずに実行を続ける。二重対が復元されても、二次側は直ちに同期状態になるわけではない。この時点で保留状態の更新を適用した後で、二次側は同期状態に戻る。一次側は、該当ボリュームに関する中止コマンドを一次DASDに出すことによって、二次側の同期を喪失させることもできる。この中止コマンドが終了し、二重対が再確立され、保留状態の更新がコピーされると、二次側は一次側と再同期する。また、オンライン保守でも、同期を喪失させることができる。

【0021】二次ボリュームが一次ボリュームと同期していない場合、二次ボリュームは、二次システムの復旧と一次側アプリケーションの再開に使用することができない。二次側の非同期ボリュームは非同期ボリュームとして識別されなければならない、二次側の復旧引継ぎ手順は、アプリケーション・アクセスを否定する（そのボリュームを強制的にオフラインにするか、そのVOLSERを変更する）ために非同期ボリュームを識別する必要がある。一次側ホストがアクセス不能になった場合に一

次側を復旧するために、二次側が呼び出されることがある。このため、二次側では、すべてのボリュームの同期状態に関するすべての関連情報が必要である。二次記憶サブシステム、すなわち、二次記憶制御装置とDASDは、一次側で検出された例外によって一次側が同期を解除する原因となるすべての条件を判別できるわけではない。たとえば、二次側が把握していない一次入出力経路またはリンクの障害のために一次側が二次側の対等機能にアクセスできない場合、一次側は二重対を解除することがある。この場合、二次側が同期状態を示すのに対し、一次側は、二重対が解除されたことを示す。

【0022】非同期二重対ボリュームが存在することを、外部通信によって二次側に通知することができる。これは、ユーザ・システム管理機能を使用することで認識できる。一次側の入出力操作はチャネル終了/装置終了/装置チェック(CE/DE/UC)状況で終了し、センス・データがエラーの特徴を示す。このような形式の入出力構成の場合、エラー回復プログラム(ERP)がエラーを処理し、入出力の完了を一次側アプリケーションに通知する前に二次プロセッサに適切なメッセージを送る。この場合、ERPの二重対中止メッセージを認識し、その情報を二次側で確保するのは、ユーザの責任である。一次側の代わりに動作可能になるよう二次側が頼りにされている場合は、始動手順によって二次DASDが二次ホストにオンライン接続され、アプリケーションの割振りのために非同期ボリュームがオンライン接続されていないことを確認するために、二次DASDサブシステムに格納された同期状況が検索される。この同期状況をすべてのERP二重対中止メッセージと組み合わせると、二次側の非同期ボリューム全体を示すピクチャが得られる。

【0023】ここで図1を参照して説明すると、同図には、一次側14と二次側15を有し、二次側15が一次側14から20キロメートル離れている、災害復旧システム10が示されている。一次側14は、そこで実行されているアプリケーションと、システム入出力およびエラー回復プログラム2(以下、入出力ERP2という)とを有するホスト・プロセッサまたは一次プロセッサ1を含む。一次プロセッサ1は、DFSMS/MVSオペレーティング・ソフトウェアを実行するIBMエンタープライズ・システム/9000(ES/9000)などでもよく、さらにそこで複数のアプリケーション・プログラムを実行することもできる。一次プロセッサ1には、IBM 3990-6型記憶制御装置などの一次記憶制御装置3がチャネル12を介して接続されている。当技術分野で既知の通り、複数のこのような一次記憶制御装置3を1台の一次プロセッサ1に接続するか、あるいは複数の一次プロセッサ1を複数の一次記憶制御装置3に接続することができる。一次記憶制御装置3には、IBM 3390 DASDなどの一次DASD4が接

続されている。複数の一次DASD4を一次記憶制御装置3に接続することができる。一次記憶制御装置3と、これに接続された一次DASD4によって、一次記憶サブシステムが形成される。また、一次記憶制御装置3と一次DASD4は、単一の一体型ユニットであってもよい。

【0024】二次側15は、チャネル13を介してIBM 3990-3型などの二次記憶制御装置6に接続されたIBM ES/9000などの二次プロセッサ5を含む。二次記憶制御装置6には、さらにDASD7が接続されている。一次プロセッサ1は、チャネル・リンクまたはT1/T3電話回線リンクなどの少なくとも1つのホスト間通信リンク11によって二次プロセッサ5に接続されている。一次プロセッサ1は、複数のエンタープライズ・システム接続(ESCON)リンク9などによって二次記憶制御装置6との直接接続を確保することもできる。その結果、入出力ERP2は、必要であれば、二次記憶制御装置6と通信可能になる。一次記憶制御装置3は、複数のESCONリンクなどの複数の対等通信リンク8を介して二次記憶制御装置6と通信する。

【0025】一次プロセッサ1で実行されるアプリケーション・プログラムによって書込み入出力操作が実行されると、入出力操作が正常に完了したことを示すハードウェア状況としてチャネル終了/装置終了(CE/DE)が出力される。入出力操作が正常に完了すると、一次プロセッサ1のオペレーティング・システム・ソフトウェアは、そのアプリケーションに書込み入出力成功のマークを付け、それにより、アプリケーション・プログラムは、最初または前の書込み入出力操作の正常終了に依存する可能性のある次の書込み入出力操作に移行できるようになる。これに対して、書込み入出力操作が不成功に終わった場合は、チャネル終了/装置終了/装置チェック(以下、CE/DE/UCという)という入出力状況が一次プロセッサ1のオペレーティング・システム・ソフトウェアに出力される。装置チェックを出力した後、入出力ERP2は、制御権を引き継ぎ、失敗した書込み入出力操作の特徴に関する具体的なセンス情報を一次記憶制御装置3から入手する。あるボリュームに固有のエラーが発生した場合は、そのエラーに関連する固有の状況が入出力ERP2に出力される。その後、入出力ERP2は、一次記憶制御装置3と二次記憶制御装置6との間、または最悪の場合は、一次プロセッサ1と二次プロセッサ5との間のデータ保全性を維持するために、新たな対等通信同期エラー回復を実行することができる。

【0026】図2および図3を参照して説明すると、同図にはエラー回復手順が示されている。図2のステップ201は、一次プロセッサ1で実行されるアプリケーション・プログラムが一次記憶制御装置3にデータ更新を送信することを含む。ステップ203では、そのデータ



更新が一次DASD4に書き込まれ、そのデータ更新が二次記憶制御装置6にシャドーイングされる。ステップ205では、二重対の状況がチェックされ、一次側と二次側が同期しているかどうかが判別される。二重対の状況が同期状態になっている場合、ステップ207でデータ更新が二次DASD7に書き込まれ、一次プロセッサ1での処理は、そこで実行されるアプリケーション・プログラムを介して続行される。

【0027】二重対が「障害発生」状態になっている場合、ステップ209で一次記憶制御装置3は、二重対で中断または障害が発生していることを一次プロセッサ1に通知する。二重対は、通信リンク8による一次記憶制御装置3と二次記憶制御装置6との通信障害によって「障害発生」状態になる場合がある。あるいは、二重対は、一次サブシステムまたは二次サブシステムいずれかのエラーによって「障害発生」状態になる場合もある。障害が通信リンク8で発生している場合、一次記憶制御装置3は、二次記憶制御装置6に直接、障害を連絡することができない。そこで、一次記憶制御装置3は、ステップ211で入出力状況としてCE/DE/UCを一次プロセッサ1に返す。入出力ERP2は、アプリケーション・プログラムを静止させ、書き込み入出力操作を要求するアプリケーションに制御権を返す前に、エラー回復とデータ保全性のためにステップ213で一次プロセッサ1の制御権を引き継ぐ。

【0028】図3は、入出力ERP2が実行する諸ステップを表している。ステップ221で入出力ERP2は一次記憶制御装置3にセンス入出力を出す。センス入出力操作は、入出力エラーの原因を記述する情報を返す。すなわち、このデータ記述情報は、具体的なエラーに関して記憶制御装置または二重対の操作に固有のものになる。一次記憶制御装置3と二次記憶制御装置6との間の対等通信リンク8で障害が発生したことがデータ記述情報によって示された場合、ステップ223で入出力ERP2は、一次記憶制御装置3および二次記憶制御装置6に対して、関係ボリュームを同期遠隔コピー障害状態に入れるように指示する記憶制御装置レベル入出力操作を出す。この二次記憶制御装置6は、複数のESCONリンク9またはホスト間通信リンク11を介して入出力ERP2から関係ボリュームの状態を受け取ることができる。その結果、二重対操作の現在の状況は、一次プロセッサ1で実行されるアプリケーションとともに、一次プロセッサ1および二次プロセッサ5の両方で維持される。コンソール18および19は、それぞれ一次プロセッサ1および二次プロセッサ5からの情報をやりとりするために設けられ、入出力ERPは、両方のコンソール18および19に状況情報を通知する。

【0029】一次記憶制御装置3および二次記憶制御装置6への同期遠隔コピー入出力操作障害が正常に完了したとき、ステップ225ではデータ保全性が維持されて

いる。このため、二次側15で復旧を試みると、二次記憶制御装置6は、「同期遠隔コピー障害」というマークを付けたボリュームを、データ回復手段（ボリューム上のそのデータの状態を判別するための従来のデータベース・ログまたはジャーナル）によってそのボリュームのデータとその同期グループ内の他のデータとの同期が取られるまで使用できないものとして識別する。

【0030】ステップ227では、同期遠隔コピー障害の状況更新について一次記憶制御装置3と二次記憶制御装置6で行われた入出力操作の正常終了を入出力ERP2が受け取ったかどうかを判別するテストが行われる。正常終了すると、入出力ERP2は、ステップ229で一次プロセッサ1に制御権を返す。正常終了していない場合は、ステップ231で次のレベルの復旧通知が行われる。この通知には、障害発生ボリュームと、一次記憶制御装置3または二次記憶制御装置6のいずれかのそのボリュームの状況が正しくない可能性があることを、コンソール18を介してオペレータに通知することが含まれる。この通知は、そこで具体的なボリューム状況を示すために、コンソール19または共用DASDデータ・セットを介して二次側15にシャドーイングされる。

【0031】ステップ233で、エラー・ログ記録データ・セットが更新される。この更新は、一次DASD4または他の記憶場所のいずれかに書き込まれ、二次側15にシャドーイングされる。このエラー回復処置が完了すると、入出力ERP2はステップ235で、書き込み入出力操作障害に関する「永続エラー」回復を一次側アプリケーションに実行させるために、一次側アプリケーションの書き込み入出力操作に「永続エラー」を通知する。エラーが修正されると、ボリューム状態は、まず保留状態（変更データの再コピー）に回復し、次に全二重に回復することができる。その後、二重対が再確立されると、データを二次DASD7に再適用することができる。

【0032】二重対を確立する場合、顧客の要求に応じて、ボリュームをCRITICALと識別することができる。CRITICALボリュームの場合、ある操作の結果、二重対の障害が発生すると、実際のエラー箇所とは無関係に、一次ボリュームの永続エラー障害が報告される。CRIT=Yの場合、障害発生対の一次DASD406に書き込もうとするその後のすべての試みは、永続エラーを受け取ることになり、対をなす二次ボリュームにシャドーイングできないデータは、その一次ボリュームに一切書き込まれなくなる。このため、必要であれば、一次側アプリケーションの処置および入出力データ操作との完全同期が可能になる。

【0033】その結果、本明細書に記載する災害復旧システム10では、入出力命令（チャンネル・コマンド・ワード（CCW））を有する一次ホスト処理エラー回復手順によって、一次および二次同期遠隔コピー・ボリュー

ムの状態を二重対から障害発生二重対へ変更できるようにする同期遠隔コピーを取り入れ、それにより、複数タイプの一次および二次サブシステム・エラーの場合にデータ保全性を維持する。アプリケーション・ベースのバックアップではなく、データ更新がリアルタイムで複写される記憶域ベースのバックアップが設けられている。また、災害復旧システム10は、(1)一次および二次記憶制御装置ボリューム状況更新、(2)オペレータ・メッセージまたはエラー・ログ記録共通データ・セットを介して具体的なボリューム更新状況に関して一次および二次ホスト・プロセッサが通知すること、および(3)CRITICALボリューム識別などの、複数レベルの一次/二次状況更新を試み、ボリューム対が障害発生二重対になる場合は、一次ボリュームへのその後の更新を防止することができる。このため、リアルタイムの完全エラー災害復旧が達成される。

#### 【0034】非同期データ・シャドーイング

非同期遠隔データ・シャドーイングは、1回の災害で一次側と二次側の両方が崩壊してしまう確率を低減するために一次側と二次側との距離をさらに大きくする必要がある場合、または一次側アプリケーションのパフォーマンスへの影響を最小限に抑える必要がある場合に使用する。一次側と二次側との距離は、現在では地球全体またはそれ以上に延長できるが、複数の一次サブシステムの背後にある複数のDASDボリュームにわたる書き込み更新を複数の二次サブシステムに同期させることは、さらに複雑である。二次記憶サブシステム上でシャドーイングするために、一次データ・ムーバを介して一次記憶制御装置から二次データ・ムーバへレコード書き込み更新を発送することができるが、両者間でやりとりされる制御データの量は、最小限でなければならない。同時に、複数の記憶制御装置に隠れている複数のDASDボリュームにわたる一次システムの場合と同様、複数の記憶制御装置にわたる二次システム上でレコード書き込み更新の順序を正確に再構築できるものでなければならない。

【0035】図4は、一次側421と遠隔側または二次側431とを含む非同期災害復旧システム400を示している。一次側421は、DFSMS/MVSホスト・ソフトウェアを実行するIBM ES/9000などの一次プロセッサ401を含む。一次プロセッサ401は、IMSおよびDBSアプリケーションなどのアプリケーション・プログラム402および403と、一次データ・ムーバ(PDM)404をさらに含む。一次プロセッサ401には、そこで実行されるすべてのアプリケーション(402、403)に共通の基準を提供するために、共通シスプレックス・クロック(sysplex clock)407が設けられ、すべてのシステム・クロックまたは時間源(図示せず)がシスプレックス・クロック407に同期し、すべての時間依存プロセスが相互に正しいタイミングで動作するようになっている。たとえば、

一次記憶制御装置405は、単一の一次記憶制御装置406への2回の連続する書き込み出力操作が同じタイム・スタンプ値を示さないように、複数のレコード書き込み更新時間を確実に区別するのに適した解像度に同期している。シスプレックス・クロック407の解像度(正確さではない)は重要である。PDM404は、シスプレックス・クロック407に接続された状態で図示されているが、書き込み出力操作がそこで発生するわけではないので、シスプレックス・クロック407に同期させる必要はない。また、一次プロセッサ401が単一の時間基準(たとえば、単一のマルチプロセッサES/9000システム)を有する場合には、シスプレックス・クロック407は不要である。

【0036】一次プロセッサ401には、IBM 3990-6型記憶制御装置などの複数の一次記憶制御装置405が光ファイバ・チャネルなどの複数のチャネルを介して接続されている。また、各一次記憶制御装置405には、IBM 3390DASDなどの複数の一次DASD406からなる少なくとも1つのストリングが接続されている。一次記憶制御装置405と一次DASD406によって、一次記憶サブシステムが形成される。各記憶制御装置405と一次DASD406は、個別のユニットである必要はなく、両者を組み合わせて単一のドロウにしてもよい。

【0037】一次側421から数千キロメートル離れた位置に配置される二次側431は、一次側421と同様に、そこで動作する二次データ・ムーバ(SDM)414を有する二次プロセッサ411を含む。あるいは、一次側と二次側が同じ場所に存在してもよく、さらに、一次データ・ムーバと二次データ・ムーバが単一のホスト・プロセッサに常駐してもよい(二次DASDは防火壁のすぐ上に設けてもよい)。二次プロセッサ411には、当技術分野で既知の通り、光ファイバ・チャネルなどのチャネルを介して複数の二次記憶制御装置415が接続されている。記憶制御装置415には、複数の二次DASD416と1つの制御情報DASD417(複数可)が接続されている。記憶制御装置415とDASD416および417によって、二次記憶サブシステムが構成される。

【0038】一次側421は、通信リンク408を介して二次側431と通信する。より具体的には、一次プロセッサ401は、仮想記憶通信アクセス方式(VTAM)通信リンク408などの通信プロトコルによって、二次プロセッサ411にデータと制御情報を転送する。この通信リンク408は、電話(T1、T3回線)、無線、無線/電話、マイクロ波、衛星などの複数の適当な通信方式によって実現できる。

【0039】非同期データ・シャドーイング・システム400は、一次DASD406へのすべてのデータ書き込みの順序が保持され、二次DASD416に適用される

(すべての一次記憶サブシステムにわたるデータ書き込み順序を保持する)ように、一次記憶制御装置405から制御データを収集する機能を含む。二次側431に送られるデータおよび制御情報は、データ保全性を保持するのに一次側421の存在が不要になるほど、十分なものでなければならない。

【0040】アプリケーション402、403は、データまたはレコード更新を生成するが、このレコード更新は、一次記憶制御装置405によって収集され、PDM404によって読み取られる。それぞれの一次記憶制御装置405は、非同期遠隔データ・シャドーイング・セッションのためにそれぞれのレコード更新をグループ化し、非特定一次DASD406のREAD要求を介してPDM404にこれらのレコード更新を提供する。一次記憶制御装置405からPDM404へのレコード更新の転送は、START入出力操作の回数および読取りから読取りの遅延を最小限にしながら、各一次記憶制御装置405と一次プロセッサ401との間で転送されるデータの量を最大にするように、PDM404によって制御され、最適化される。PDM404は、非特定READ間の時間間隔を変えることで、一次記憶制御装置とホストとのこの最適化だけでなく、二次DASD416用のレコード更新の通用期間も制御することができる。

【0041】データ保全性を維持しながら、PDM404がレコード更新を収集し、そのレコード更新をSDM414に送信するには、すべての一次記憶サブシステムにおいて二次DASD416に対して行われる一次DASD406のレコードWRITEシーケンスを再構築するのに十分な制御データとともに、特定の期間の間、適切な複数の時間間隔でレコード更新を送信する必要がある。一次DASD406のレコードWRITEシーケンスの再構築は、自己記述レコードをPDM404からSDM414に渡すことによって達成される。SDM414は、所与の時間間隔分のレコードが紛失しているかどうか、または不完全になっているかどうかを判別するために、その自己記述レコードを検査する。

【0042】図5および図6は、接頭部ヘッダ500(図5)と、一次記憶制御装置405によって生成されたレコード・セット情報600(図6)とを含む、各自己記述レコードごとにPDM404が作成するジャーナル・レコード形式を示している。各自己記述レコードは、それぞれの時間間隔の時間順に二次DASD416に適用できるように、それぞれの時間間隔ごとにさらにSDM414によってジャーナル処理される。

【0043】ここで図5を参照して説明すると、各レコード・セットの先頭に挿入される接頭部ヘッダ500は、接頭部ヘッダ500と、各レコード・セットごとにSDM414に送信される実際の一次レコード・セット情報600との長さの合計を記述するための総データ長501を含む。操作タイム・スタンプ502は、PDM

404が現在処理している操作セットの開始時間を示すタイム・スタンプである。この操作タイム・スタンプ502は、1組の一次記憶制御装置405に対してREAD RECORD SET機能を実行する際に(シスプレックス・クロック407に応じて)PDM404によって生成される。一次DASD406の書き込みの入出力時間610(図6)は、各一次記憶制御装置405のREAD RECORD SETごとに固有のものである。操作タイム・スタンプ502は、すべての記憶制御装置で共通のものである。

【0044】READ RECORD SETコマンドは、PDM404によって出されるが、以下の条件のいずれかの場合に予測できる。

- (1) 一次記憶制御装置405の所定のしきい値に基づく、その一次記憶制御装置のアテンション割込み
- (2) 所定の時間間隔に基づく、一次プロセッサ401のタイマ割込み
- (3) レコード・セット情報が、使用可能であるがまだ読み取られていない未解決のレコード・セットに関する追加情報を示す場合

条件(2)では、タイマ間隔を使用して、低レベル活動の期間中に二次システム431がどの程度遅れて実行するかを制御する。条件(3)は、PDM404が一次記憶制御装置405の活動に遅れないようにするためにさらに活動を駆動する処理間隔中に、PDM404がすべてのレコード・セットを待ち行列処理しなかった場合に発生する。

【0045】時間間隔グループ番号503は、現行レコード・セット(整合性グループのうちの所与の時間間隔グループについてすべての一次記憶制御装置405にわたるレコードのセット)が属す時間間隔(操作タイム・スタンプ502とレコード読取り時間507によって境界が示される)を識別するためにPDM404が出力する。グループ内順序番号504は、所与の時間間隔グループ503内の各レコード・セットごとに一次記憶制御装置405用のアプリケーションWRITE入出力の順序を(PDM404に対して)識別するためにハードウェアが提供するIDに基づいて導出される。一次SSID(補助記憶域ID)505は、各レコード・セットごとに一次記憶制御装置405の特定の一次記憶制御装置を明確に識別するものである。二次ターゲット・ボリューム506は、パフォーマンス上の考慮事項に応じて、PDM404またはSDM414のいずれかによって割り当てられる。レコード読取り時間507は、すべての一次記憶制御装置405に共通の操作タイム・スタンプを提供し、現行間隔のレコード・セットの終了時間を示す。

【0046】操作タイム・スタンプ502およびレコード読取り時間507は、各一次記憶制御装置405から得た複数組の読取りレコード・セットをグループ分けす

るためにPDM404が使用する。複数組の読取りレコード・セットをグループ分けするための同期はPDM404にだけ重要であるため、PDM404は、シスプレックス・クロック407に接続されていないPDM404だけを動作させる中央処理装置（CPU）クロックに同期してもよい。PDM404はレコード更新を書き込まないが、前述の通り、レコード更新は共通の時間源に同期していなければならない。

【0047】次に図6を参照して説明すると、レコード・セット情報600は、一次記憶制御装置405によって生成され、PDM404によって収集される。更新固有情報601～610は、レコード更新が行われた実際の一次DASD406を含む各レコードの一次装置ユニット・アドレス601を含む。シリンダ番号／ヘッド番号（CCHH）602は、各レコード更新ごとの一次DASD406上の位置を示す。一次記憶制御装置のセッションIDである一次SSID603は、一次SSID505と同じものである。状況フラグ604は、特定のデータ・レコード620が後に続くかどうかに関する状況情報を提供する。順序番号605および630は、レコード・セット全体（PDM404に転送されたすべてのデータ）が読み取られたかどうかを示すために各レコードに番号を1つずつ割り当てる。一次DASD書き込み出力タイプ606は、各レコードについて行われた書き込み操作のタイプを識別する操作標識である。この操作標識は、更新書き込み、フォーマット書き込み、部分トラック・レコード・フォロー、完全トラック・データ・フォロー、消去コマンド実行、または全書き込み実行を含む。検索指数607は、最初に読み取られたレコード・セット・データ・レコード620に関する初期位置決め情報を示す。セクタ番号608は、レコードが更新されたセクタを識別する。カウント・フィールド609は、後続の特定のレコード・データ・フィールド620の数を記述する。一次DASD406の書き込み更新が行われたホスト・アプリケーション時間は、更新時間610に記録される。特定のレコード・データ620は、各レコード更新ごとのカウント／キー／データ（CKD）フィールドを提供する。最後に、順序番号630は、読み取られたレコード・セット全体がPDM404に転送されたかどうかを示すために順序番号605と比較される。

【0048】一次DASD406でレコード更新が書き込まれたのと同じ順序でSDM414がそのレコード更新をコピーできるように、ソフトウェア・グループが呼び出した整合性グループで更新レコードが処理される。整合性グループを作成するのに使用する情報（すべての記憶制御装置405から収集したすべてのレコード・セットにわたる）は、操作タイム・スタンプ502、時間間隔グループ番号503、グループ内順序番号504、一次制御装置SSID505、レコード読取り時間507、一次装置アドレス601、一次SSID603、お

よび状況フラグ604を含む。1つの時間間隔グループ用のすべてのレコードがSDM414側で各記憶制御装置405ごとに受け取られたかどうかを判別するのに使用する情報は、時間間隔グループ番号503、グループ内順序番号504、物理制御装置ID505、および一次SSID603および各操作時間間隔ごとに各一次記憶制御装置405から返される読取りレコード・セットの総数を含む。完全復旧可能な一次DASD406のレコード更新と同等に二次DASD416上にレコード更新を配置するのに必要な情報は、二次ターゲット・ボリューム506、CCHH602、一次DASD書き込み入出力タイプ606、検索指数607、セクタ番号608、カウント609、更新時間610、および特定のレコード・データ620を含む。

【0049】図7および図8は、復旧時間とジャーナル転送時間を単純化した現行ジャーナル内容を記述するための状態テーブル700とマスタ・ジャーナル800をそれぞれ示す。状態テーブル700は、PDM404とSDM414が収集し、両者に共通の構成情報を提供し、一次記憶制御装置のセッションID（SSID番号）およびその制御装置でのボリュームと、対応する二次記憶制御装置のセッションIDおよび対応するボリュームとを含む。このため、構成情報は、どの一次ボリューム710または一次DASDエクステントが二次ボリューム711または二次DASDエクステントにマッピングされるかを追跡する。状態テーブル700まで単純に拡張して部分ボリューム・エクステント712（CCHHからCCHHまで）を示す場合、部分ボリューム遠隔コピーは、ここに記載するのと同じ非同期遠隔コピー方法を使用して達成できるが、完全ボリュームの場合より細分性（トラックまたはエクステント）はより細くなる。

【0050】マスタ・ジャーナル800は、整合性グループ番号、ジャーナル・ボリューム上の位置、および操作タイム・スタンプを含む。また、マスタ・ジャーナル800は、整合性グループにグループ化した特定のレコード更新を維持する。状態テーブル700とマスタ・ジャーナル800は、災害復旧をサポートするため、一次システム410がもはや存在しないスタンドアロン環境で動作できなければならない。

【0051】制御項目全体が正しく書き込まれるようにするため、タイム・スタンプ制御は各マスタ・ジャーナル800の前後に置かれる。このタイム・スタンプ制御は、さらに二次DASD417に書き込まれる。制御要素は二重項目（1）および（2）を含み、次に示す例のように一方の項目が必ず現行項目になる。

（1）タイム・スタンプ制御 | 制御情報 | タイム・スタンプ制御

（2）タイム・スタンプ制御 | 制御情報 | タイム・スタンプ制御

いかなる時点でも、(1)または(2)のいずれかの項目が現行または有効項目になるが、有効項目は前後に等しいタイム・スタンプ制御を持つ項目である。災害復旧では、制御情報を得るために、最新のタイム・スタンプを持つ有効項目を使用する。この制御情報は、状態情報（記憶制御装置、装置、および適用される整合性グループに関する環境情報）とともに、二次記憶制御装置415にどのレコード更新が適用されたかを判別するのに使用する。

#### 【0052】整合性グループ

所定の時間間隔の間にすべての一次記憶制御装置405にわたるすべての読取りレコード・セットが二次側431で受け取られると、SDM414は、受け取った制御情報を解釈し、レコード更新が最初に一次DASD406上で書き込まれたのと同じ順序でそのレコード更新が適用されるように、受け取った読取りレコード・セットをレコード更新・グループとして二次DASD416に適用する。このため、一次側アプリケーションの順序（データ保全性）整合性はすべて二次側431で維持される。このプロセスは、以下、整合性グループの形成と呼ぶ。整合性グループの形成は、次のような仮定に基づいて行われる。(A) 独立しているアプリケーション書込みが制御装置の順序命令に違反しない場合は、どのような順序でもアプリケーション書込みを実行できる。

(B) 従属しているアプリケーション書込みは、タイム・スタンプの順に実行しなければならないため、アプリケーションは、書込み番号1から制御装置終了、装置終了を受け取る前に従属書込み番号2を実行することができない。(C) 第二の書込みは必ず(1)遅いタイム・スタンプを持つ第一の書込みと同じレコード・セット整合性グループに入るか、(2)後続のレコード・セット整合性グループに入る。

【0053】図9を参照して説明すると、同図には、記憶制御装置SSID1、SSID2、およびSSID3など（記憶制御装置はいくつでも含めることができるが、この例では明解にするため3つ使用する）に関する整合性グループの形成例（整合性グループは一次側421または二次側431のいずれにも形成できるはずである）が示されている。時間間隔T1、T2、およびT3は昇順に発生するものと想定する。時間間隔T1の操作タイム・スタンプ502は、記憶制御装置SSID1、SSID2、およびSSID3について設定されている。PDM404は、時間間隔T1～T3の間に記憶制御装置SSID1、2、および3からレコード・セット・データを入手する。時間間隔T1のSSID1、2、および3に関するレコード・セットは、時間間隔グループ1であるG1（時間間隔グループ番号503）に割り当てられる。グループ内順序番号504は、SSID1、2、および3のそれぞれについて示され、この場合、SSID1は11:59、12:00、および1

2:01に3つの更新を持ち、SSID2は12:00および12:02に2つの更新を持ち、SSID3は11:58、11:59、および12:02に3つの更新を持つ。時間間隔T2およびT3のレコード・セットは列挙されているが、簡略化のため、更新時間の例は示されていない。

【0054】ここで、二次側431で受け取った制御情報およびレコード更新に基づいて、整合性グループNを生成することができる。時間間隔グループ番号1のレコード更新が時間間隔グループ番号2のレコード更新より遅くならないようにするため、記憶制御装置SSID1、2、および3のそれぞれの最後のレコード更新の最も早い読取りレコード・セット時間と等しい、最小時間が設定される。この例では、最小時間は12:01になる。最小時間と等しいかそれ以上の読取りレコード・セット時間を有するレコード更新はすべて整合性グループN+1に含まれる。1つのボリュームに対する2つのレコード更新時間が等しい場合、シスプレックス・クロック407の十分な解像度が与えられる可能性はほとんどないが、時間間隔グループN内の早い順序番号を持つレコード更新は、整合性グループN用のそのグループとともに保管される。ここで、レコード更新は、読取りレコード・セット時間に基づいて順序づけされる。複数のレコード更新の時間が等しい場合、小さい順序番号を持つレコード更新は、大きい順序番号を持つレコード更新の前に置かれる。これに対して、複数のレコード更新のタイム・スタンプが等しいが、ボリュームが異なる場合は、そのレコード更新が同じ整合性グループに保管されている限り、任意の順序にすることができる。

【0055】一次記憶制御装置405が、指定の時間間隔の間に読取りレコード・セットへの応答を完了しなかった場合、その一次記憶制御装置405が完了するまで、整合性グループを形成することはできない。一次記憶制御装置405がその操作を完了しなかった場合は、未着割込みのために、システムの未着割込みハンドラが制御権を受け取り、操作が終了する。これに対して、一次記憶制御装置405が適切な時間に操作を完了した場合は、入出力が完了に至り、通常操作が続行される。整合性グループの形成では、一次記憶制御装置405に対する書込み操作にタイム・スタンプが付けられると予想される。しかし、プログラムによっては、タイム・スタンプが付けられずに書込みが生成されるものもある。この場合、一次記憶制御装置405は、タイム・スタンプとしてゼロを返す。整合性グループの形成は、データが読み取られたタイム・スタンプに基づいて、タイム・スタンプを持たないこれらのレコードの境界を示すことができる。整合性グループの時間別にレコード更新の境界を容易に示せないほど、タイム・スタンプを持たないレコード更新が一定の時間間隔の間に多数発生した場合、二重ボリュームが同期していないというエラーが発生す

る可能性がある。

【0056】図10および図11は、整合性グループを形成する方法を示す流れ図である。図10を参照して説明すると、このプロセスは、ステップ1000から始まり、一次側421が、行うべき遠隔データ・シャドーイングを確立する。ステップ1010では、シスプレックス・クロック407を同期クロック（図4）として使用して、すべてのアプリケーション入出力操作にタイム・スタンプが付けられる。PDM404は、ステップ1020で各一次記憶制御装置405との遠隔データ・シャドーイング・セッションを開始するが、このステップは、データまたはレコードがシャドーイングされる一次ボリュームの識別を含む。ステップ1030では、各アプリケーションWRITE入出力操作（図6を参照）ごとに一次記憶制御装置405によってレコード・セット情報600がトラッピングされる。

【0057】ステップ1040は、前述の通り、アテンション・メッセージを含むプロンプト、所定のタイミグ間隔、または読取りレコード数増加の通知に応じて、PDM404が、捕捉したレコード・セット情報600を各一次記憶制御装置405から読み取ることを含む。ステップ1050でPDM404がレコード・セットの読取りを開始すると、PDM404は、各レコード・セットの前に特定のジャーナル・レコード（ジャーナル・レコードは、接頭部ヘッダ500と、レコード・セット情報600を含む）を作成するための接頭部ヘッダ500（図5を参照）を付ける。このジャーナル・レコードには、二次側431（または一次側421）で整合性グループを形成するのに必要な制御情報（およびレコード）が含まれる。

【0058】ステップ1060では、PDM404が通信リンク408を介して（整合性グループがそこで形成される場合は、同じデータ・ムーバ・システム内で）SDM414に生成したジャーナル・レコードを送信する。SDM414は、ステップ1070で状態テーブル700を使用し、データ・シャドーイング・セッション用に確立した各時間間隔グループおよび一次記憶制御装置405ごとに、受け取ったレコード更新をグループ番号別および順序番号別に収集する。ステップ1080でSDM414は、ジャーナル・レコードを検査し、各時間間隔グループごとにすべてのレコード情報を受け取ったかどうかを判別する。ジャーナル・レコードが不完全な場合は、ステップ1085によって、SDM414はPDM404に必要なレコード・セットを再送信するよう通知する。PDM404が正しく再送信できない場合は、二重ボリューム対に障害が発生している。ジャーナル・レコードが完全な場合は、SDM414による整合性グループの形成を含むステップ1090が実行される。

【0059】図11を参照すると、同図には、整合性グ

ループを形成するためのステップ1090（図10）を表すステップ1100～1160が示されている。整合性グループの形成は、ステップ1100から始まるが、このステップでは、各ソフトウェア整合性グループが二次DASD417（図4）上のSDM414ジャーナル・ログ（"hardened"）に書き込まれる。ステップ1110は、時間間隔グループが完全かどうかを判別するテストを実行する。すなわち、各一時記憶制御装置405は、少なくとも1つの読取りレコード・セット・バッファを提示したか、レコード・セット・バッファ内にこのようなレコード更新が置かれていないという確認をPDM404から受け取らなければならない、しかもデータ（またはヌル）を持つすべての読取りレコード・セット・バッファがSDM414によって受け取られていなければならない。時間間隔グループが不完全な場合は、ステップ1110は、必要なデータが受け取られるまで、一次記憶制御装置405からのレコード・セットの読取りを再試行する。エラーが発生した場合は、特定の1つまたは複数の二重ボリューム対に障害が発生している可能性がある。完全な時間間隔グループを受け取ると、ステップ1120は、第一の整合性グループ・ジャーナル・レコードを判別する。この第一（または現行）の整合性グループ・ジャーナル・レコードとは、最も早い操作タイム・スタンプ502と、同じ操作タイム・スタンプ502を持つすべてのレコードの最も早い更新時間610を含むレコードである。

【0060】ステップ1130は、現行整合性グループ・ジャーナル・レコードに含まれるレコードを検査して、どのレコードがそこに最後に含めるレコードかを判別する（一部のレコードは除去され、次の整合性グループ・ジャーナル・レコードに含まれる）。現行整合性グループ・ジャーナル・レコードの最後のレコードは、各一次記憶制御装置405ごとに最大更新時間のうちの最小更新時間（最小時間）として判別される（つまり、各一次記憶制御装置405の最後の更新が比較され、これらのうちの最も早いものだけが現行整合性グループ・ジャーナル・レコードに残る）。

【0061】現行整合性グループ・ジャーナル・レコードに残っているこれらのレコード更新は、ステップ1140で、更新時間610とグループ内順序番号504に応じて順序付けされる。レコード更新を持たない一次記憶制御装置405は、整合性グループに関与しない。ステップ1150では、現行整合性グループに残っているレコード更新（最小時間より遅い更新時間を持つもの）が、次の整合性グループに渡される。それぞれのグループ内順序番号504は、空バッファで終わり、その操作時間間隔の間にすべての読取りレコード・セットが読み取られたことを示すはずである。空バッファがない場合は、現行ソフトウェア整合性グループ内の最後のレコードを定義するステップ1120を、レコード読取り時間

507および更新時間610と併せ使用して、一次記憶制御装置405におけるアプリケーションWRITE入出力操作の正しい順序を決定することができる。

【0062】ステップ1160は、完全災害復旧の制約の下で特定の書き込み更新が二次DASD416に適用される、遠隔データ・シャドーイング・プロセスのバックエンドを表している。二次DASD416に更新内容を書き込む際に入出力エラーが発生するか、または二次側431全体が停止し、最初期設定された場合は、書き込みプロセスに入っていた整合性グループ全体を最初から再適用することができる。このため、どの二次DASD416の入出力が行われたか、どの入出力が行われなかったか、どの入出力が処理中かなどを追跡せずに、遠隔シャドーイングを行うことができる。

【0063】二次入出力書き込み  
ステップ1160の重要な構成要素は、二次側431が一次側421から後れをとらないように、PDM414によってレコードが二次DASD416に効率よく書き込まれることである。必要な効率率は、主に、様々な二次DASD416への複数の入出力操作を同時に実行することによって達成される。二次DASD416を一度に1つずつ連続して書き込むと、二次側431は一次側421からかなり遅れてしまう恐れがある。単一チャネル・コマンド・ワード(CCW)連鎖を介して単一の二次装置宛ての整合性グループごとにレコードを書き込めば、二次側431ではさらに高い効率を得られる。それぞれの単一CCW連鎖内では、そこで行われる各二次DASD416のデータ・トラックへの入出力操作が一次ボリュームでの発生順に維持されている限り、その入出力操作をさらに最適化することができる。

【0064】特定の整合性グループ用の二次入出力操作を単一CCW連鎖内で最適化する場合、一部は一次書き込み入出力操作のパターンに基づいて行われ、一部は二次DASD416の物理特性に基づいて行われる。最適化は、二次DASD416がカウント/キー/データ(CKD)か、拡張カウント/キー/データ(ECKD)か、固定ブロック方式(FBA)かなどに応じて、多少変化する可能性がある。その結果、所与の時間間隔の間に1つの一次DASD406に対して行われる複数のWRITE入出力(m)は、1つの二次DASD416のボリュームに対する単一のSTART入出力操作に削減することができる。このように二次記憶制御装置415に対するSTART入出力の回数をm:1に最適化すると、二次DASD416は後れをとらずに済み、それにより、一次側421のレコード更新をもっと精密にシャドーイングすることができる。

【0065】正常な遠隔データ・シャドーイングとそれによる二次入出力の最適化の重要点は、一貫性のあるコピーを復旧に使用できるように、二次DASD416に対して同時に行う複数の入出力操作のいずれかで発生す

る回復不能エラーを最小限にすることである。所与の二次書き込みで失敗すると、その後の従属書き込みが条件付け書き込みを伴わずに記録される恐れがある(たとえば、実際にはデータベース用の実際の更新書き込みが失敗に終わっているのにデータベース・レコードが更新されたことを示すログ項目は、二次DASD416のコピーの順序整合性に違反する)。

【0066】その更新失敗が復旧されるまで、失敗した二次DASD416のコピーはアプリケーションの復旧に使用できなくなる。失敗した更新は、SDM414によってPDM404から現行コピーを要求することで修正できるはずである。その間に二次データ・コピーは不整合になり、そのため、PDM404が現行更新で応答し、それ以前の他のすべての更新がPDM404によって処理されるまで使用できなくなる。通常、失敗した更新の復旧に要する時間は、十分な災害復旧保護のために受け入れられないほど長い非復旧ウィンドウを示す。

【0067】効果的な二次側431の入出力最適化は、所与の整合性グループについて書き込まれるデータ・レコード・セットを検査し、ECKD方式などの二次DASD416の特定の方式の規則に基づいて連鎖を構築することで実現される。ここに開示する最適化技法は、入出力エラーが発生した場合に整合性グループを適用する際に、CCW連鎖を再実行できるように、または、二次初期プログラム・ロード(IPL)復旧の場合に、データを紛失せずに整合性グループ全体を再適用できるように、入出力エラーからの復旧を単純化するものである。

【0068】図12は、ECKD方式用のすべてのWRITE入出力の組合せに対応するCCW連鎖を構築するための完全整合性グループ復旧(FCGR)の規則を要約して示すもので、ここではCCHHRレコード形式が使用される(シリンダ番号、ヘッド番号、レコード番号)。図12は、1つの整合性グループの範囲内でDASDトラックに対して行われるWRITE入出力操作の可能な組合せをそれぞれ検査することによって作成される。図12のFCGR規則(図14および図15に記載する)は、整合性グループを適用する際のエラーについて完全復旧を行うためにデータ配置(二次DASD416の入出力書き込みCCW連鎖)を管理する場合に従うものである。図12に示すFCGR規則は、新しいWRITE入出力操作が追加されるたびに適切に拡張されることになる。これらの規則は、二次側431のハードウェアまたはソフトウェアで実現することができる。FCGR規則は、同一DASDトラック分析に対するREADレコード・セットを、一次DASD406のWRITE入出力タイプ、検索指数、およびカウント・フィールドとキー・フィールドの検査に還元するので好都合である。

【0069】図12に示すように整合性グループの書き込み操作を検査せずにDASDトラックが書き込まれる

と、前に書き込まれたデータ・レコードの再書き込みができなくなる可能性がある。たとえば、以下の内容の連鎖があると想定する。

レコード5へのWRITE UPDATE

レコード1へのFORMAT WRITE

この場合、レコード1とレコード5は同じDASDトラックに存在し、レコード1がレコード5の前に置かれる。レコード5はUPDATE WRITE CCWによって更新されるが、FORMAT WRITE入出力CCWは、トラックの残りを消去してレコード1を更新するため、レコード5は削除される。この連鎖を再実行しなければならない場合、レコード5の先頭に配置するLOCATERECORD CCWが位置決めポイントを持たなくなる（レコード5が存在しなくなる）ため、この連鎖は先頭から完全に回復することができない。一次側421ですでに書き込み操作が正常に行われているので、データの整合性と保全性を維持するには、二次DASD416上の整合性グループ全体をいつでも適用することが必要である。

【0070】図16のステップ1410～1470は、図11のステップ1160によって表され、図12に定義されるFCGR規則を使用するプロセスの詳細を示すものである。ステップ1410でSDM414は、現行整合性グループの各種レコードを2つのカテゴリに分割する。第一のカテゴリは、同じ二次DASDボリューム向けの入出力命令を含み、第二のカテゴリは、第一のカテゴリに含まれるレコードのうち、同じCCHH向けのレコードの入出力命令を含む（すなわち、同じDASDトラックに更新されるレコード）。

【0071】現行整合性グループのレコードのカテゴリ分類後、ステップ1420では、トラック上のデータ配置を識別し、トラック／レコードのアドレス指定を行うために、アプリケーションのWRITE入出力およびSDM414のWRITE入出力を二次DASD416の方式、たとえば、ECKD方式のFCGR規則（図12を参照）に適合させる。SDM414は、ステップ1430で、同じボリュームに対する二次DASDのWRITE入出力操作を単一出力CCW連鎖にグループ分けする。ステップ1440は、実際の二次DASD416の書き込み用の検索指数と特定のレコード・データ（CKDフィールド）に応じて、それぞれの二次DASD416のヘッド・ディスク・アセンブリ（HDA）を移動させることを含む。

【0072】ステップ1450では、後続の書き込み操作によって、前の書き込み操作またはDASD検索指数（ここで消去されるレコードでの分割など）を無効にするかどうかを判別するために図12のFCGR規則を使用して、第二のカテゴリ（通常、レコードを受け取るトラックごとに1つずつ、複数の第二のカテゴリが存在する）を構成するこれらのレコードについてREAD SET

BUFFERS1と2を比較する。READ SET BUFFERS1と2は隣接する読取りレコード・セットが含まれている。FCGR規則に従うと、エラーが発生した場合に、一次側421からレコード更新を再度受け取らなくても、SDM414は整合性グループ全体を再書き込みできるようになる。SDM414が現行整合性グループを二次DASD416に適用した後、ステップ1460では、状態テーブル（図7）とマスタ・ジャーナル（図8）を更新する。

【0073】ステップ1470は、次の整合性グループ（現行整合性グループになるもの）を獲得して、処理をステップ1410に戻すので、遠隔コピー・プロセスはリアルタイムで続行される。一次側421から二次側431への通信が終了した場合、遠隔コピー・プロセスは停止する。この通信は、ボリューム対がPDM404によってプロセスから削除された場合、一次側が破壊された場合（災害が発生した場合）、規則的な運転停止が行われた場合、または二次側431で特定の引継ぎ処置が行われた場合に終了することがある。二次側431でジャーナル処理された整合性グループは、引継ぎ操作中に二次DASD416に適用することができる。一次側421が捕捉したデータのうち、SDM414によって完全に受け取られていないデータだけが紛失する。

【0074】要約すると、これまで同期および非同期遠隔データ二重化システムについて説明してきた。非同期遠隔データ二重化システムは、記憶域ベースのリアルタイム・データ・シャドーイングを提供する。一次側は、レコード更新を生成するアプリケーションを実行し、一次側から離れた位置にある二次側は、レコード更新をシャドーイングして、一次側の災害復旧を行う。この非同期遠隔データ二重化システムは、一次側の時間依存プロセスを同期させるためのシスプレックス・クロックと、アプリケーションを実行するための一次側の一次プロセッサを含み、一次プロセッサはそこに一次データ・ムーバを有している。一次プロセッサには、各レコード更新ごとに書き込み入出力操作を出すために複数の一次記憶制御装置が連結され、それぞれの一次記憶制御装置DASD書き込み入出力操作がシスプレックス・クロックに同期している。複数の一次記憶装置はこの書き込み入出力操作を受け取り、それに応じてレコード更新をそこに格納する。一次データ・ムーバは、各レコード更新ごとに複数の一次記憶制御装置からレコード・セット情報を収集し、所定のグループのレコード・セット情報に接頭部ヘッダを付加する。この接頭部ヘッダと所定のグループのレコード・セット情報が自己記述レコード・セットを形成する。それぞれのレコード・セット情報は、一次装置アドレス、シリンダ番号とヘッド番号（CCHH）、レコード更新順序番号、書き込み入出力タイプ、検索指数、セクタ番号、およびレコード更新時間を含む。接頭部ヘッダは、総データ長、操作タイム・スタンプ、時間間隔



グループ番号、およびレコード読取り時間を含む。二次側の二次プロセッサは二次データ・ムーバを有し、この二次データ・ムーバが一次側から自己記述レコード・セットを受け取る。二次プロセッサには複数の二次記憶制御装置が連結され、二次記憶制御装置にはレコード・更新のコピーを格納するために複数の二次記憶装置が連結されている。二次データ・ムーバは、送られてきた自己記述レコード・セットが完全なものであるかどうかを判定し、その自己記述レコード・セットから整合性グループを形成し、さらに、レコード更新が複数の一次記憶装置に書き込まれた順序と整合する順序で複数の二次記憶装置に書き込むために、各整合性グループから得たレコード更新を複数の二次記憶制御装置に出力する。

【0075】特に本発明の実施例に関連して本発明を図示し説明してきたが、当業者は、本発明の精神および範囲を逸脱せずに形態および詳細の様々な変更が可能であることに留意されたい。たとえば、整合性グループは、受け取った自己記述レコード・セットに基づいて二次データ・ムーバによって形成されたものとして説明してきたが、書き込みレコード・セットに基づいて一次側で整合性グループを形成するか、二次側の他の部分で形成することもできる。一次側と二次側の記憶装置の形式は、同じである必要はない。たとえば、CKDレコードを固定ブロック方式(FBA)タイプのレコードなどに変換することも可能である。また、記憶装置は、DASD装置に限定されているわけではない。

【0076】

【発明の効果】本発明の実施により、災害復旧のためにDASDデータを二次側にシャドーイングするための改良された設計および方法を提供することができる。

【0077】まとめとして、本発明の構成に関して以下の事項を開示する。

【0078】(1) 一次プロセッサで実行される1つまたは複数のアプリケーションによって生成されたデータ更新が一次記憶サブシステムによって受け取られ、一次記憶サブシステムは入出力書き込み操作によって各データ更新を書き込まれ、各書き込み入出力操作にタイム・スタンプが付けられ、共通タイマによってタイム・スタンプの同期が取られ、一次プロセッサと同じ地域にあるか一次プロセッサから離れた位置にあるかにかかわらず、二次システムは、災害復旧のために二次側が使用できるように順序整合性のある順序でデータ更新をシャドーイングする。災害復旧機能を提供するために整合性グループを形成する方法において、前記方法が、(a) 一次記憶サブシステムで発生する各書き込み入出力操作にタイム・スタンプを付けるステップと、(b) 各データ更新ごとに一次記憶サブシステムから書き込み入出力操作レコード・セット情報を収集するステップと、(c) データ更新とそれぞれのレコード・セット情報から自己記述レコード・セットを生成し、この自己記述レコード・セット

が、二次システムだけによる書き込み入出力操作の順序を再生成するために十分な制御情報を含んでいるステップと、(d) 自己記述レコード・セットを間隔グループにグループ分けし、各間隔グループが、操作タイム・スタンプ開始時間から測定され、所定の間隔しきい値の間持続するステップと、(e) 最も早い操作タイム・スタンプを有する自己記述レコード・セットの間隔グループとして現行整合性グループを選択し、一次記憶サブシステムでの入出力書き込み操作の時間順に基づいて、個々のデータ更新が現行整合性グループ内で順序付けされるステップを含む方法。

(2) 各間隔グループの開始時間を識別するための操作タイム・スタンプに基づいて、一次記憶サブシステムとのセッションを開始し、各間隔グループの境界が連続する操作タイム・スタンプによって示されることが、ステップ(b)にさらに含まれることを特徴とする、上記(1)に記載の方法。

(3) 各間隔グループを記述する接頭部ヘッダを追加することがステップ(d)に含まれることを特徴とする、上記(1)に記載の方法。

(4) 自己記述レコード・セットからなる間隔グループを二次側に送信するステップ(f)をさらに含むことを特徴とする、上記(3)に記載の方法。

(5) 受け取った各自己記述レコード・セットが完全なものであるかどうかを二次側で判定するステップ(g)をさらに含むことを特徴とする、上記(4)に記載の方法。

(6) 自己記述レコード・セットが不完全であると二次側が判定した場合に、欠落データ更新を再送信するよう二次側が一次側に要求することが、ステップ(g)にさらに含まれることを特徴とする、上記(5)に記載の方法。

(7) 各時間間隔グループが完全なものであるかどうかを二次側で判定するステップ(h)をさらに含むことを特徴とする、上記(6)に記載の方法。

(8) 間隔グループが不完全であると二次側が判定した場合に、欠落レコード・セットを再送信するよう二次側が一次側に要求することが、ステップ(h)にさらに含まれることを特徴とする、上記(7)に記載の方法。

(9) 整合性グループで順序づけされたように、対応する一次側の書き込み入出力操作の順序に応じて、二次側で受け取ったデータ更新を二次記憶サブシステムに書き込むステップ(i)をさらに含むことを特徴とする、上記(8)に記載の方法。

(10) 災害復旧のために遠隔データ・シャドーイングを行うシステムにおいて、このシステムは、一次データ・ムーバと、レコード更新を生成するアプリケーションとを実行する一次プロセッサを有する一次側を含み、一次プロセッサは、一次プロセッサから一次記憶サブシステムに出される書き込み入出力操作に応じてレコード更新

を格納するための記憶装置を有する一次記憶サブシステムに連結され、一次側は、一次側の時間依存操作を同期させるために共通のシステム・タイマをさらに含み、システムは、一次プロセッサと通信する二次プロセッサと、順序整合性のある順序でレコード更新のコピーを格納するための二次記憶サブシステムとを有する二次側をさらに含み、(a)一次記憶サブシステムの各書込み入出力操作にタイム・スタンプを付けるステップと、

(b)一次記憶サブシステム内の各記憶装置とのセッションを確立するステップと、(c)一次記憶サブシステム内の各記憶装置からレコード・セット情報を収集するステップと、(d)レコード・セットとそれぞれのレコード・セット情報を一次データ・ムーバに読み込むステップと、(e)各レコード・セットの前にヘッダを付け、それに基づく自己記述レコード・セットを生成するステップと、(f)所定の時間間隔に応じて、時間間隔グループ単位で自己記述レコード・セットを二次プロセッサに送信するステップと、(g)自己記述レコード・セットから整合性グループを形成するステップと、

(h)順序整合性のある順序で各整合性グループのレコード更新を二次記憶サブシステムにシャドーイングするステップとを含む、順序整合性のある順序でレコード更新をシャドーイングするための方法。

(11)レコード・セットが非同期的に二次プロセッサに送信されることを特徴とする、上記(10)に記載の方法。

(12)ステップ(g)が二次側で行われることを特徴とする、上記(10)に記載の方法。

(13)受け取った各自己記述レコード・セットが完全であるかどうかを二次側で判定するステップが、ステップ(f)にさらに含まれることを特徴とする、上記(10)に記載の方法。

(14)受け取った自己記述レコード・セットが不完全であると二次側が判定した場合に、欠落レコード更新を再送信するよう一次側に要求するステップが、ステップ(f)にさらに含まれることを特徴とする、上記(13)に記載の方法。

(15)各時間間隔グループが完全であるかどうかを二次側で判定するステップ(i)をさらに含むことを特徴とする、上記(10)に記載の方法。

(16)時間間隔グループが不完全であると二次側が判定した場合に、欠落レコード・セットを再送信するよう一次側に要求することが、ステップ(i)にさらに含まれることを特徴とする、上記(15)に記載の方法。

(17)ステップ(c)が、レコード・セット情報において、各レコード更新が格納されている一次記憶装置上の物理的位置を識別することを特徴とする、上記(10)に記載の方法。

(18)ステップ(c)が、レコード・セット情報において、セッション内に一次記憶装置に格納された各レコ

ード更新の順序と更新時間を識別することを特徴とする、上記(17)に記載の方法。

(19)ステップ(e)が、接頭部ヘッダにおいて、セッション用の間隔グループ番号と、そこで参照される各レコード更新用のグループ内順序を識別することを特徴とする、上記(10)に記載の方法。

(20)1つまたは複数のアプリケーションを実行する一次プロセッサを有し、この1つまたは複数のアプリケーションがレコード更新を生成し、一次プロセッサがそれに基づく自己記述レコード・セットを生成し、自己記述レコード・セットが二次システムに送られ、二次システムがリアルタイム災害復旧のために自己記述レコード・セットに基づいて順序整合性のある順序でレコード更新をシャドーイングし、一次プロセッサが一次記憶サブシステムに連結され、一次記憶サブシステムがレコード更新を受け取って、そこに各レコード更新を格納するために書込み入出力操作を実行し、一次プロセッサが、同期を取るためにアプリケーションと一次記憶サブシステムに共通の時間源を提供するためのシスプレックス・クロックと、各レコード更新ごとにレコード・セット情報を提供するように一次記憶サブシステムに指示し、複数のレコード更新と、それに対応するそれぞれのレコード・セット情報を時間間隔グループにグループ分けし、それに接頭部ヘッダを挿入し、各時間間隔グループがレコード・セットを自己記述する一次データ・ムーバ手段とを含む一次システム。

(21)一次記憶サブシステムが、書込み入出力操作を出す複数の一次記憶制御装置と、複数の一次記憶サブシステムに連結された複数の一次記憶装置とを含むことを特徴とする、上記(20)に記載の一次システム。

(22)複数の一次記憶装置が直接アクセス記憶装置であることを特徴とする、上記(21)に記載の一次システム。

(23)一次データ・ムーバ手段が、各時間間隔グループに関与する複数の一次記憶制御装置のうちの各一次記憶制御装置用のそれぞれの書込み入出力操作ごとにレコード・セット情報を収集することを特徴とする、上記(22)に記載の一次システム。

(24)一次プロセッサにおいて、各書込み入出力操作にシスプレックス・クロックに関するタイム・スタンプが付けられ、各書込み入出力操作が、複数の一次記憶制御装置のうちの1つの一次記憶制御装置に出され、各一次記憶制御装置が、タイム・スタンプを保持し、対応する読取りレコード・セットに入れてそのタイム・スタンプを一次データ・ムーバ手段に返すことを特徴とする、上記(23)に記載の一次システム。

(25)各レコード・セット情報が、複数の一次記憶装置のうちの1つの一次記憶装置上における対応するレコード更新の物理的位置を識別することを特徴とする、上記(24)に記載の一次システム。

(26) 各レコード・セット情報が、対応するレコード更新の一次サブシステムID、一次装置アドレス、シリンダ番号、およびヘッド番号を識別することを特徴とする、上記(24)に記載の一次システム。

(27) 一次データ・ムーバ手段が、1つの時間間隔グループに参与するすべての一次記憶制御装置にわたる各書込み入出力更新の相対順序を識別することを特徴とする、上記(24)に記載の一次システム。

(28) レコード更新をジャーナル処理し、一次システムおよび二次システム上の各レコード更新の記憶場所を相互参照するために、一次データ・ムーバ手段が状態テーブルを作成することを特徴とする、上記(27)に記載の一次システム。

(29) 一次データ・ムーバ手段が、二次システムに状態テーブルを送信することを特徴とする、上記(27)に記載の一次システム。

(30) 一次側と二次側を含み、二次側が災害復旧のために一次側のレコード更新をリアルタイムでシャドーイングし、レコード更新が一次側で実行されるアプリケーションによって生成され、一次側が、シスプレックス・クロックと、レコード更新を生成するアプリケーションを実行し、各レコード更新ごとに対応する書込み入出力操作を出し、一次データ・ムーバをそこに有する一次プロセッサと、レコード更新を格納するよう指示され、各レコード更新ごとに出された書込み入出力操作を実行する複数の一次記憶制御装置と、対応する書込み入出力操作に応じて、レコード更新を受け取ってそこに格納する複数の一次記憶装置とを含み、書込み入出力操作が互いに正しい順序で並べられるように、シスプレックス・クロックによって同期が取られた通りに、一次プロセッサによって一次プロセッサと各書込み入出力にタイム・スタンプが付けられ、一次データ・ムーバが、複数組のレコード更新を収集し、複数の一次記憶制御装置のうちのそれぞれによって提供された各レコード・セット情報と対応するレコード更新を組み合わせ、各レコード・セット情報が、それぞれの対応する書込み入出力操作の相対順序と時間を含み、一次データ・ムーバが時間間隔グループ別にレコード更新を収集して、各時間間隔グループに接頭部ヘッダを挿入し、接頭部ヘッダが各時間間隔グループに含まれるレコード更新を識別する情報を含み、各レコード・セット情報と接頭部ヘッダが、自己記述レコード・セットを生成するために組み合わせられ、自己記述レコード・セットが二次側に送信され、自己記述レコード・セットが、一次側からの追加の通信がなくても順序整合性のある順序でレコード更新をそこにシャドーイングするために十分な情報を二次側に提供する、遠隔データ・シャドーイング・システム。

(31) 遠隔データ・シャドーイングに参与するために識別されたすべての一次記憶制御装置とのセッションを確立することによって、一次データ・ムーバが時間間隔

グループを形成することを特徴とする、上記(30)に記載の遠隔データ・シャドーイング・システム。

(32) 一次データ・ムーバが、識別されたすべての一次記憶制御装置からレコード・セット情報を収集することを特徴とする、上記(31)に記載の遠隔データ・シャドーイング・システム。

(33) 一次プロセッサが、二次側に自己記述レコードを送信することを特徴とする、上記(32)に記載の遠隔データ・シャドーイング・システム。

(34) 一次データ・ムーバが、自己記述レコードから整合性グループを形成することを特徴とする、上記(32)に記載の遠隔データ・シャドーイング・システム。

(35) 二次側が、送信された自己記述レコードから整合性グループを形成することを特徴とする、上記(33)に記載の遠隔データ・シャドーイング・システム。

(36) レコード更新をジャーナル処理し、一次システムおよび二次システム上の各レコード更新の記憶場所を相互参照する状態テーブルを、一次データ・ムーバが作成することを特徴とする、上記(32)に記載の遠隔データ・シャドーイング・システム。

(37) 複数の一次記憶装置が直接アクセス記憶装置(DASD)であることを特徴とする、上記(32)に記載の遠隔データ・シャドーイング・システム。

(38) 各レコード・セット情報が、一次装置アドレスと、シリンダ番号およびヘッド番号(CCHH)と、レコード更新順序番号と、書込み入出力タイプと、検索指数と、セクタ番号と、レコード更新時間とを含むことを特徴とする、上記(37)に記載の遠隔データ・シャドーイング・システム。

(39) 接頭部ヘッダが、総データ長と、操作タイム・スタンプと、時間間隔グループ番号と、レコード読取り時間とを含むことを特徴とする、上記(37)に記載の遠隔データ・シャドーイング・システム。

(40) レコード更新を生成するアプリケーションを実行する一次側を含み、一次側から離れた位置に二次側を有し、二次側がレコード更新をシャドーイングして、一次側に災害復旧を提供する、記憶域ベースのリアルタイム・データ・シャドーイングを行う非同期遠隔データ二重化システムにおいて、非同期遠隔データ二重化システムが、一次側の時間依存プロセスを同期させるためのシスプレックス・クロックと、アプリケーションを実行し、対応するレコード更新用の書込み入出力操作を出し、一次データ・ムーバをそこに有する、一次側の一次プロセッサと、各レコード更新ごとに書込み入出力操作を1つずつ受け取り、それぞれの一次記憶制御装置書込み入出力操作が一次プロセッサによってシスプレックス・クロックと同期している、複数の一次記憶制御装置と、対応する書込み入出力操作に応じて、レコード更新をそこに格納するための複数の一次記憶装置とを含み、一次データ・ムーバが、各レコード更新ごとに複数の一

次記憶制御装置からレコード・セット情報を収集して、  
所定のグループのレコード・セット情報に接頭部ヘッダ  
を付加し、接頭部ヘッダと所定のレコード・セット情報  
グループが自己記述レコード・セットを形成し、各レコ  
ード・セット情報が、一次装置アドレス、シリンダ番号  
およびヘッド番号（CCHH）、レコード更新順序番  
号、書き込み出力タイプ、検索指数、セクタ番号、およ  
びレコード更新番号を含み、接頭部ヘッダが、総データ  
長、操作タイム・スタンプ、時間間隔グループ番号、およ  
びレコード読取り時間を含み、二次データ・ムーバを  
有し、その二次データ・ムーバが一次側から自己記述レ  
コード・セットを受け取る、二次側の二次プロセッサ  
と、二次プロセッサに連結された複数の二次記憶制御装  
置と、レコード更新を格納する複数の二次記憶装置とを  
さらに含み、二次データ・ムーバが、送信された自己記  
述レコード・セットが完全なものであるかどうかを判定  
し、自己記述レコード・セットから整合性グループを形  
成し、複数の一次記憶装置にレコード更新が書き込まれ  
たときの順序に整合する順序で複数の二次記憶装置に書  
き込むために各整合性グループから得たレコード更新を  
複数の二次記憶制御装置に出力する、非同期遠隔データ  
二重化システム。

【図面の簡単な説明】

【図1】同期遠隔データ・シャドーイング機能を有する  
災害復旧システムのブロック図である。

【図2】図1の災害復旧システムにより同期遠隔コピー  
を提供する方法を示す流れ図である。

【図3】入出力エラー回復プログラム（入出力ERP）  
操作の方法を示す流れ図である。

【図4】非同期遠隔データ・シャドーイング機能を有す  
る災害復旧システムのブロック図である。

【図5】図4の一次側からの読取りレコード・セットの  
前に付く接頭部ヘッダを示すデータ形式図である。

【図6】読取りレコード・セットを構成する各種フィー  
ルドを示すデータ形式図である。

【図7】ボリューム構成情報を識別する状態テーブルで  
ある。

【図8】図4の二次側が使用するマスタ・ジャーナルで

ある。

【図9】整合性グループを形成するためのシーケンス例  
である。

【図10】整合性グループを形成するために情報および  
読取りレコード・セットを収集する方法を示す流れ図で  
ある。

【図11】整合性グループを形成する方法を示す流れ図  
である。

【図12】DASDトラックに対する所与の入出力操作  
シーケンスの場合のECKD方式装置用の完全整合性グ  
ループ復旧規則アプリケーションを示すテーブルであ  
る。

【図13】図12のテーブルで使用する規則の説明の構  
成を示す図である。

【図14】図12のテーブルで使用する規則の説明の一  
部である。

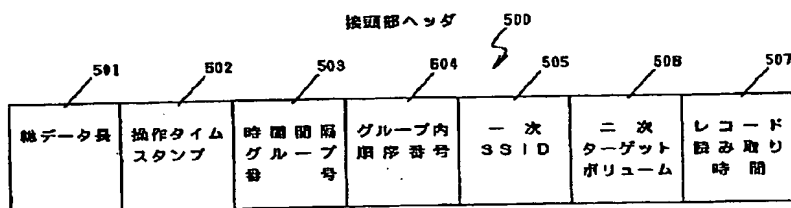
【図15】図12のテーブルで使用する規則の説明の一  
部である。

【図16】完全整合性グループ復旧機能を持つ二次側に  
読取りレコード・セット・コピーを書き込む方法の流れ  
図である。

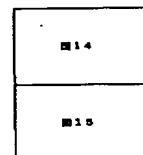
【符号の説明】

- 1 一次プロセッサ
- 2 入出力ERP
- 3 一次記憶制御装置
- 4 一次DASD
- 5 二次プロセッサ
- 6 二次記憶制御装置
- 7 二次DASD
- 8 対等通信リンク
- 9 エンタープライズ・システム接続（ESCON）リ  
ンク
- 10 災害復旧システム
- 11 ホスト間通信リンク
- 12 チャンネル
- 13 チャンネル
- 14 一次側
- 15 二次側

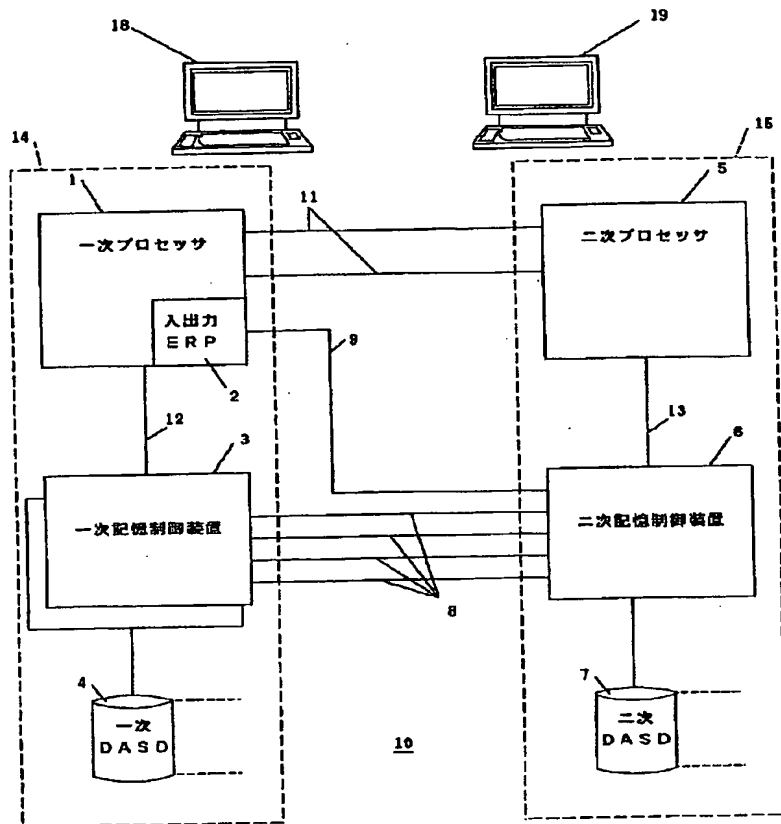
【図5】



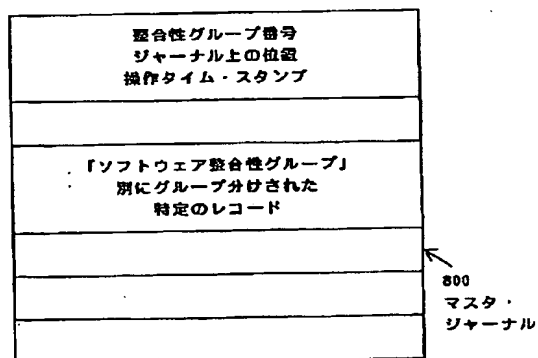
【図13】



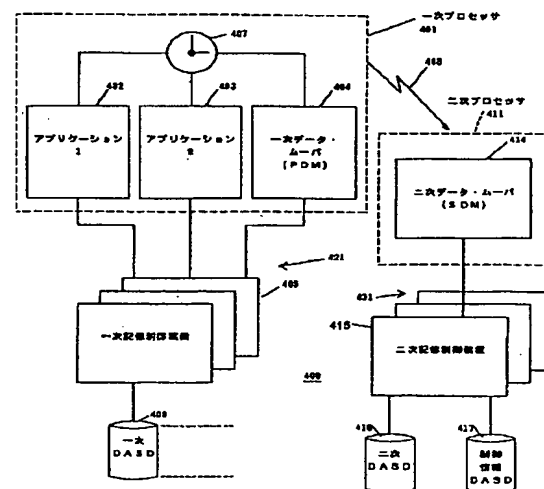
【図1】



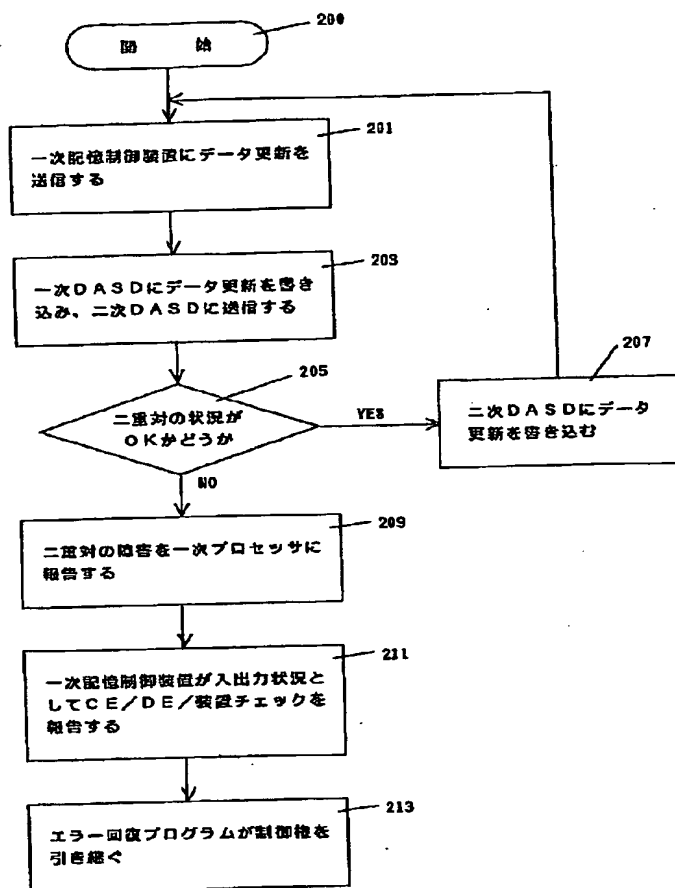
【図8】



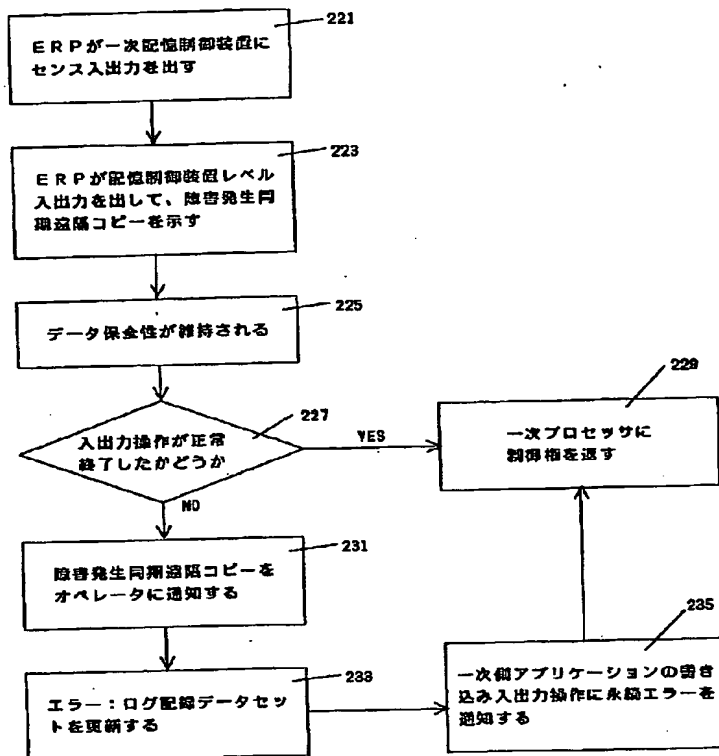
【図4】



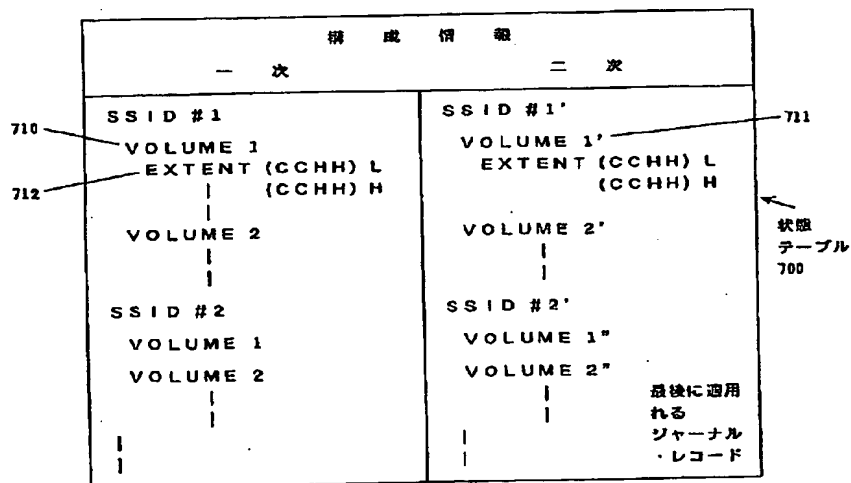
【図2】



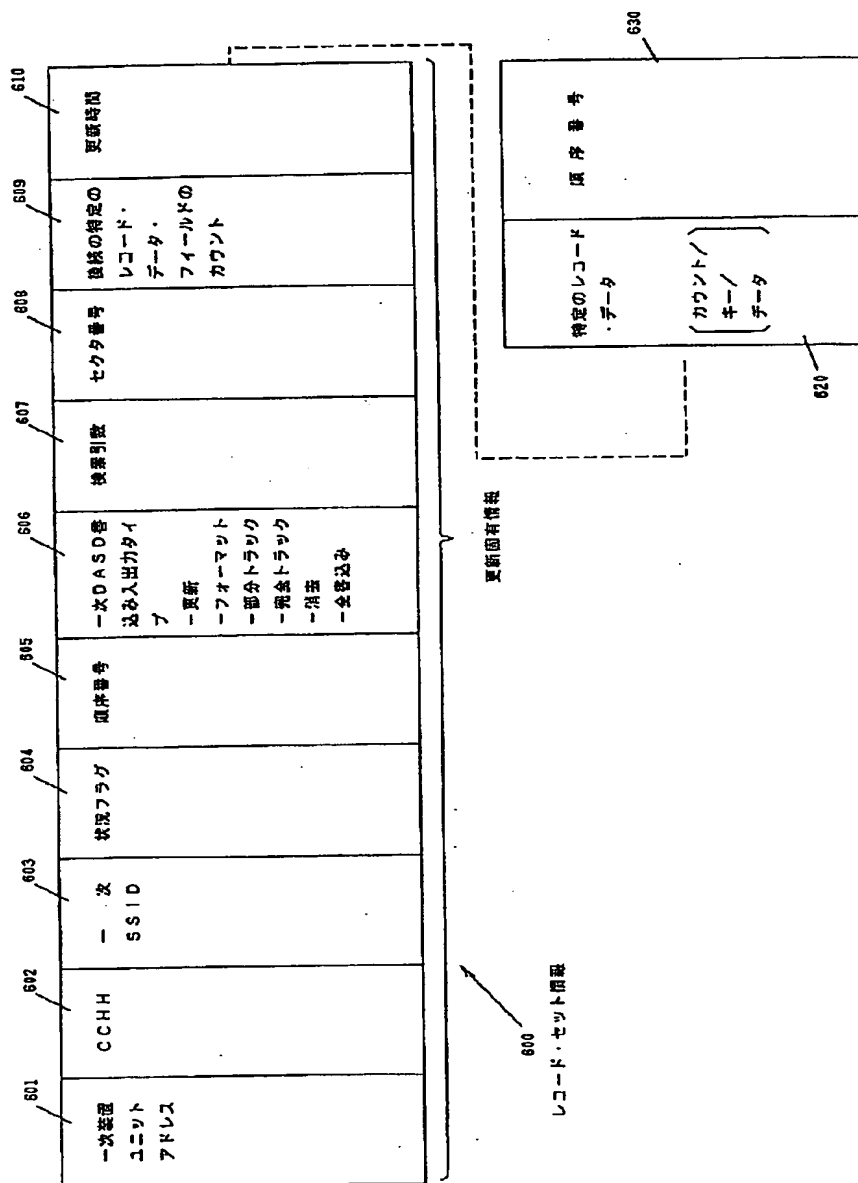
【図3】



【図7】



【図6】

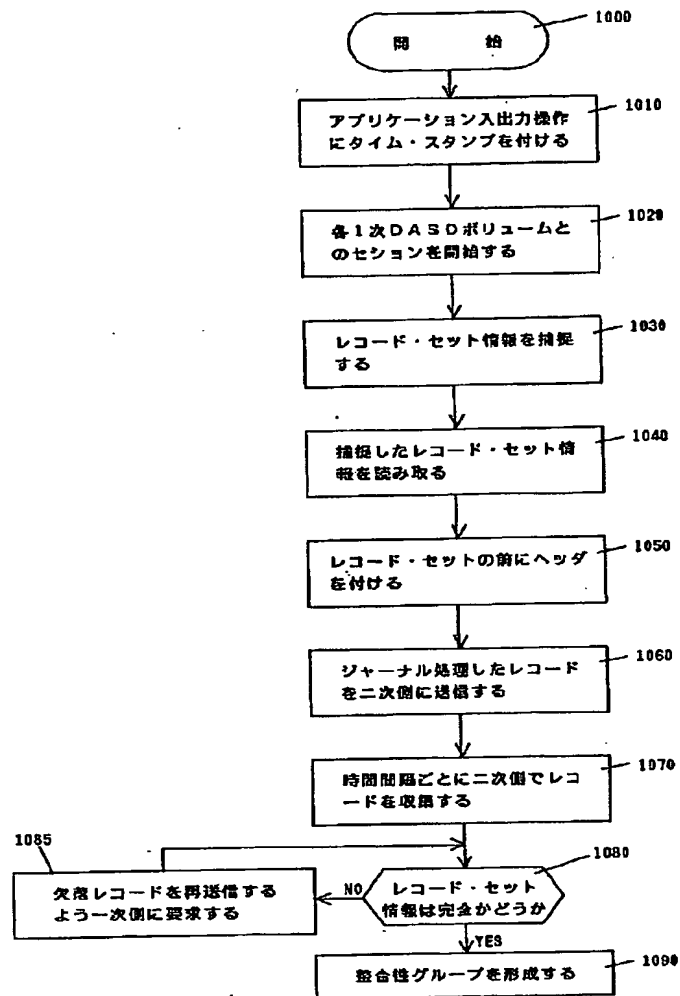




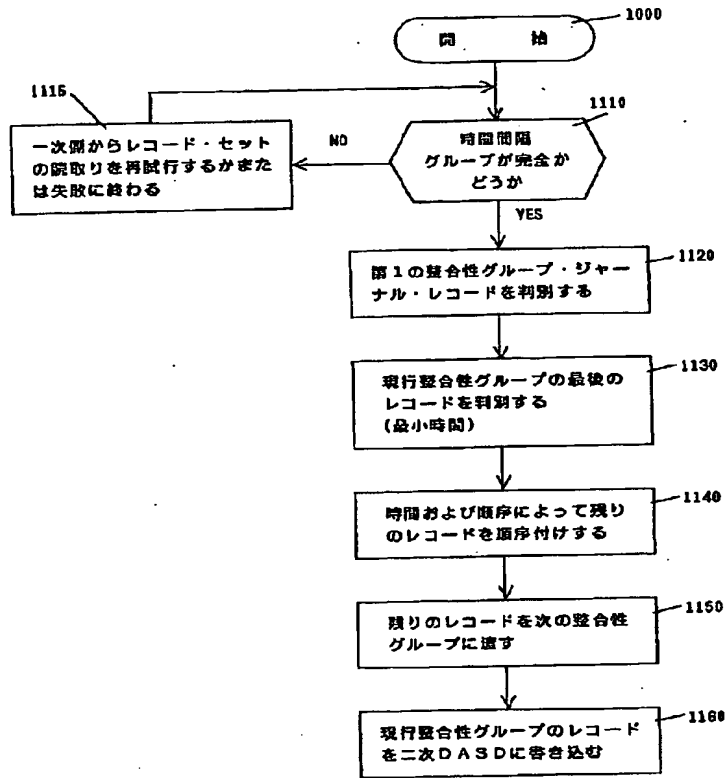
【図9】

| 物理制御<br>装置ID                             | 操作タイム<br>・スタンプ | 時間間隔<br>グループ<br>番号 | 読取りレコード・セット<br>更新時間/制御装置 |         |         |
|--|----------------|--------------------|--------------------------|---------|---------|
|  |                |                    | 順序3の1                    | 順序3の2   | 順序3の3   |
| SSID1                                    | T1             | G1                 | 11:59 ②                  | 12:00 ⑤ | 12:01 ⑥ |
| SSID2                                    | T1             | G1                 | 12:00 ④                  | 12:02 ⑦ |         |
| SSID3                                    | T1             | G1                 | 11:58 ①                  | 11:59 ③ | 12:02 ⑧ |
| SSID1                                    | T2             | G2                 |                          |         |         |
| SSID2                                    | T2             | G2                 |                          |         |         |
| SSID3                                    | T2             | G2                 |                          |         |         |
| SSID3                                    | T3             | G3                 |                          |         |         |
| 整合性グループ番号1                               |                |                    |                          |         |         |
| ① 11:58                                  |                |                    |                          |         |         |
| ② 11:59                                  |                |                    |                          |         |         |
| ③ 11:59                                  |                |                    |                          |         |         |
| ④ 12:00                                  |                |                    |                          |         |         |
| ⑤ 12:00                                  |                |                    |                          |         |         |
| ⑥ 12:01                                  |                |                    |                          |         |         |
| 読取り順に並べた場合の最も早い操作時間T1<br>SSID全体で最も早い更新時間 |                |                    |                          |         |         |
| SSID全体での最大更新時間の最小時間                      |                |                    |                          |         |         |
| 整合性グループ番号2                               |                |                    |                          |         |         |
| ⑦ -----                                  |                |                    |                          |         |         |
| ⑧ -----                                  |                |                    |                          |         |         |

【図10】



【図11】



【図12】

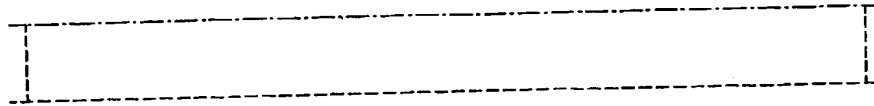
完全整合性グループ回復規則

| 既取りレコード・セット・バックアップ番号2 |                 | 既取りレコード・セット・バックアップ番号1 |                  |                  |      |      |                |                |  |
|-----------------------|-----------------|-----------------------|------------------|------------------|------|------|----------------|----------------|--|
| 入出力書き込み操作のタイプ         |                 |                       |                  |                  |      |      |                |                |  |
|                       | 書き込み更新・<br>KL=0 | 書き込み更新・<br>KL≠0       | 完全フォーマット<br>書き込み | 部分フォーマット<br>書き込み | 完全消去 | 部分消去 | 任意書き込み<br>KL=0 | 任意書き込み<br>KL≠0 |  |
| 書き込み更新<br>KL=0        | W <sup>*</sup>  | E <sup>*</sup>        | N                | J                | D    | K    | W              | E <sup>*</sup> |  |
| 書き込み更新<br>KL≠0        | E <sup>*</sup>  | W <sup>*</sup>        | N                | J                | D    | K    | E <sup>*</sup> | W              |  |
| 完全フォーマット書き込み          | T               | T                     | R                | T                | R    | T    | R              | R              |  |
| 部分フォーマット書き込み          | C               | C                     | N                | H                | P    | L    | W              | W              |  |
| 完全消去                  | T               | T                     | R                | R                | R    | T    | T              | T              |  |
| 部分消去                  | B               | B                     | N                | M                | E    | G    | W              | W              |  |
| 任意書き込み<br>KL=0        | W               | E <sup>*</sup>        | W                | W                | E    | W    | W              | E <sup>*</sup> |  |
| 任意書き込み<br>KL≠0        | E <sup>*</sup>  | W                     | W                | W                | E    | W    | E <sup>*</sup> | W              |  |

【図14】

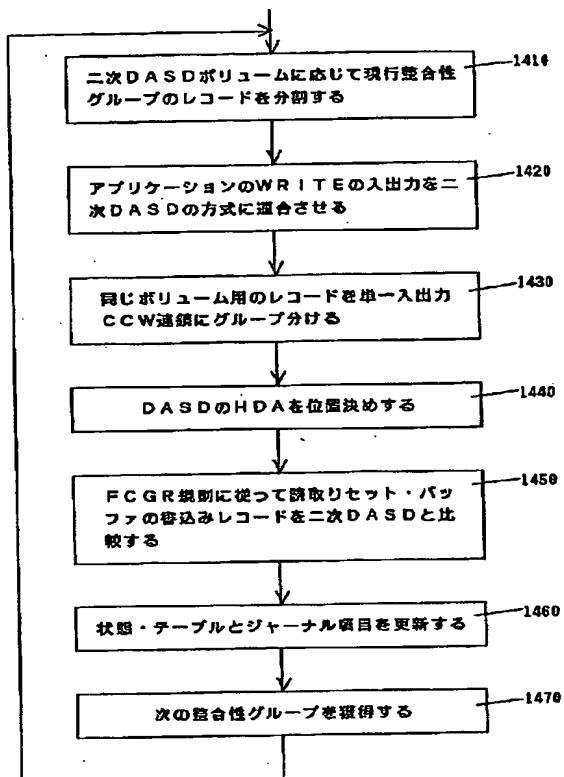
- B - 番号1の検索が番号2の検索回数より大きい場合、番号1を捨てるが両方とも実行する
  - C - 番号1のレコードが番号2に第1のレコードと等しいがそれ以上である場合、番号1を捨てる  
が両方とも実行する
  - D - 番号2がR0を更新している場合、両方とも実行するがエラーになる
  - E - エラー（発生してはならない）
  - E' - 番号1と番号2が同じレコードである場合、エラーになる  
（両者間で形式読み込みを行わずに発生してはならない）
  - F - 番号2の第1のレコードがR1である場合、両方とも読み込むが、エラーになる
  - G - 番号1の検索回数が番号2の検索回数と等しいかそれ以上である場合、番号1を捨てるがエラー  
になる
  - H - 番号1の検索回数が番号2の最後のレコードより大きい場合、番号1を捨てる。または番号2の  
検索回数が番号1の最後のレコードより大きい場合、エラーになるが両方とも読み込む
- 
- X - さらに最適化するには以下の手順を実行できる  
（番号1の検索回数が番号2の最後のレコードと等しいかそれ以上である）または（番号1の最  
後のレコードが番号2の最後のレコードと等しいかそれ以上である）しかも（番号2の検索回  
数が番号1の最後のレコードと等しいか、それ以下である）場合番号1を捨てる  
または（番号2の検索回数が番号1の最後のレコードより大きい）  
場合  
エラーになる  
または両方とも読み込む

【図15】



- J - 番号2のレコード（または検索指数）が番号1の最後のレコードより大きい場合、エラーになるか両方とも書き込む
- K - 番号2のレコード（または検索指数）が番号1の検索指数より大きい場合、エラーになるか両方とも書き込む
- L - 番号1の検索指数が番号2の検索指数と等しいかそれ以上である場合、両方とも書き込むかまたは番号1を捨てるかあるいはエラーになる
- M - (番号1の検索指数が番号2の検索指数と等しいかそれ以上である) 場合、番号1を捨てる  
または両方とも書き込む
- N - 番号2の検索指数が番号1の最後のレコードより大きい場合、エラーになるか両方とも書き込む
- R - 番号1を捨ててもよい
- T - 番号1を捨てなければならない
- W - 両方とも書き込む
- W\* - 番号1と番号2が同じレコードを持つ場合、番号1を捨てるかまたは両方とも実行するかまたはレコードを組み合わせて一方の書き込みを実行する

【図16】



フロントページの続き

|         |  |    |   |  |
|---------|--|----|---|--|
| (72)発明者 | ロナルド・マイナード・カーン<br>アメリカ合衆国85748 アリゾナ州ツー<br>ソン ノース・コレット・プレイス<br>761                  | 35 | (56)参考文献  | 特開 平7-6099 (JP, A)<br>特開 平5-204739 (JP, A)<br>特開 平4-33027 (JP, A)<br>特開 昭55-87262 (JP, A)<br>特開 昭64-67675 (JP, A)<br>特開 昭63-138441 (JP, A)<br>特開 平5-334161 (JP, A) |
| (72)発明者 | グレゴリー・エドワード・マックブライ<br>ド<br>アメリカ合衆国85715 アリゾナ州ツー<br>ソン イースト・フェアマウント・プレ<br>イス 8622   | 40 |   |  |
| (72)発明者 | デヴィッド・マイケル・シャクルフォ<br>ード<br>アメリカ合衆国85705 アリゾナ州ツー<br>ソン ウェスト・サーバー・プレイス・<br>ドライブ 1348 | 45 | (58)調査した分野(Int. Cl. <sup>7</sup> , DB名)<br>G06F 3/06<br>G06F 11/16 - 11/20<br>G06F 12/00, 12/16 |  |